

Б.И. КРЫЛОВ,

В.Д. БОБКОВ, Д.И. МОНАСТЫРНЫЙ

**Вычислительные
методы
высшей
математики**



В. И. КРЫЛОВ, В. В. БОБКОВ, П. И. МОНАСТЫРНЫЙ

Вычислительные методы высшей математики

ТОМ

1

Допущено
Министерством высшего
и среднего специального
образования БССР в качестве
учебного пособия для факуль-
тетов прикладной математики
университетов

Издательство «Вышэйшая школа». Минск 1972

518

K85

УДК 518.12 (075.8)

Рецензенты:

кафедра вычислительной математики математико-механического факультета Ленинградского университета (зав. кафедрой докт. физ.-мат. наук проф. *М. К. Гавурич*); акад. *А. Н. Тихонов*

Научный редактор

докт. физ.-мат. наук проф. *И. П. Мысовских*

Крылов В. И. и др.

K85 Вычислительные методы высшей математики. Т. 1. Под ред. И. П. Мысовских. Мн., «Вышэйш. школа», 1972.

584 с. с илл.

Книга является первым томом учебного пособия по теории вычислительных методов математики для университетов. Она будет полезна также для студентов технических учебных заведений с достаточно большой программой математики. Вместе с тем книга рассчитана на широкий круг лиц, интересующихся теорией методов вычислений.

2-2-4
8-71

518

Крылов Владимир Иванович, Бобков Владимир Васильевич, Монастырный Петр Ильич

Под редакцией

Мысовских Ивана Петровича

ВЫЧИСЛИТЕЛЬНЫЕ МЕТОДЫ ВЫСШЕЙ МАТЕМАТИКИ. ТОМ 1

Редактор *Т. Майборода*. Худож. редактор *В. Валентович*. Техн. редактор *М. Кислякова*.

Корректоры *А. Белянкина, В. Козлова*.

АТ 04238. Сдано в набор 13/VII 1971 г. Подписано к печати 30/XI 1971 г. Бумага 70×90¹/₁₆ типогр. № 1. Печ. л. 36,5 (42,705). Уч.-изд. л. 38,11. Изд. № 70—61. Зак. 367. Тираж 10 000 экз. Цена 1 руб. 47 коп. Издательство «Вышэйшая школа» Государственного комитета Совета Министров БССР по печати. Редакция литературы по естествознанию и математике. Минск, ул. Кирова, 24.

Ордена Трудового Красного Знамени типография Издательства ЦК КП Белоруссии. Минск, Ленинский пр., 79.

ПРЕДИСЛОВИЕ

Авторы стремились написать учебное пособие по теории вычислительных методов математики, предназначенное для университетов и доступное для студентов технических учебных заведений с широкой программой математики.

В основу пособия положены лекции, читавшиеся авторами для студентов специальности вычислительной математики в Ленинградском и Белорусском государственных университетах. Основной курс лекций содержал лишь наименьший объем сведений, обязательный для всех студентов этой специальности. Авторы считали необходимым дополнить его некоторыми вопросами теории вычислительных методов, которые позволили бы более полно изложить отдельные разделы теории и довести их до вида, в какой-то мере приближающегося к современному состоянию их в науке.

Мы включили в книгу также отдельные вопросы, которые, по нашему убеждению, полезно знать тем студентам, кто будет заниматься в будущей своей работе подготовкой научных и технических задач к численному их решению.

Все такие дополнительные вопросы излагались авторами в специальных курсах.

Наконец, для изложения отдельных тем требовались сведения из анализа и алгебры, которые не всегда входят в программы обязательных курсов математики. Авторы стояли перед выбором: нужно было либо дать краткое изложение недостающих сведений в этой книге, либо отсылать к специальным книгам и журнальным статьям, что в большинстве случаев весьма затруднило бы читателя, так как пришлось бы собирать нужный материал нередко по кускам на большом числе страниц.

Авторы выбрали первую из этих возможностей и предпочли недостающие сведения поместить в книгу, изложив их по возможности в кратком виде. В тех случаях, когда эти сведения можно было органически связать с основными вопросами, они вносились в соответствующие тексты. Это

оказалось возможным сделать в небольшом числе случаев. Примером могут служить дополнительные сведения из линейной алгебры, которые читатель найдет в начале гл. 2. Когда же дополнения являлись инородным телом в тексте, авторы выносили их в конец книги в форме добавлений к основному тексту.

Основной текст разделен на две неравные части. К первой из них отнесено все, принадлежащее основному курсу теории вычислительных методов. Эта часть набрана обычным шрифтом. Вторая часть текста, набранная петитом, содержит дополнительные вопросы, о которых только что говорилось.

Мы считаем, что весь объем сведений, который мы хотели бы включить в пособие, удастся разместить в двух томах приблизительно одинаковых объемов. Подробное содержание первого тома указано в оглавлении, второй же том будет посвящен изложению вычислительных методов решения дифференциальных уравнений, как обыкновенных, так и с частными производными, а также методов решения интегральных уравнений, теории улучшения сходимости рядов и последовательностей и изложению некоторых вопросов построения общей теории вычислительных методов на основе функционального анализа. Вторым том выйдет в свет в 1973 г.

Авторы хотели бы сделать некоторые замечания о характере изложения. Пособие предназначено для лиц, приступающих к изучению вычислительных методов и ранее не знакомых с ними. Мы считали поэтому необходимым начать изложение каждого из методов с подробного описания идеи, на которой этот метод основан, и условий, при которых от него можно ожидать удовлетворительной точности результатов. Там, где это было можно сделать, мы стремились выяснить наглядным путем характер изменения погрешности метода в зависимости от числа шагов, величины шага или других параметров метода. Все это позволяло выяснить те черты метода, которые образуют его качественную характеристику. И только после этого мы переходили к изложению теорем, выясняющих условия сходимости метода или устанавливающих оценку его погрешности. Такие теоремы дают преимущественно более глубоко лежащую количественную характеристику. Обе эти характеристики мы считали одинаково важными и каждой из них старались уделить достаточное место в изложении.

В пособии нет численных примеров. Отказаться от них нас побудили следующие соображения. Если стремиться не только изложить теорию вычислительных методов, но и научить студентов их применению к решению задач, необходимо ввести в книгу достаточно большое число примеров с подробным объяснением как способов выбора методов вычисления, так и техники вычислений. А это сильно увеличило бы объем и привело бы к механическому объединению теоретического пособия с руководством для практикума.

Лицам, которые будут самостоятельно заниматься изучением вычислительных методов по нашей книге, авторы хотят сообщить некоторые

сведения об уровне знаний, на которые рассчитано изложение каждого раздела. Прежде всего различные главы книги будут требовать при чтении различных знаний. Кроме того, основной текст книги, напечатанный обычным шрифтом, потребует от читателя сравнительно небольшого запаса знаний, дополнительный же текст, набранный петитом, рассчитан на более высокий уровень знаний.

Ниже приводятся сведения об уровне необходимых знаний как по главам, так и по частям текста.

Для чтения основного текста гл. 1 достаточно знания университетского курса анализа в объеме трех семестров или курса математики высшего технического учебного заведения с широкой программой математического образования.

Чтение петита потребует дополнительного знания элементов теории метрических пространств и теории операторов. Авторы рекомендуют перед чтением петита просмотреть добавление I к книге, где можно найти большую часть нужных сведений.

Для чтения гл. 2 и 3 достаточно знать основные теоремы о системах линейных уравнений, матрицах, их собственных значениях и векторах, которые содержатся в университетских программах алгебры первых трех семестров.

Все необходимые дополнительные вопросы изложены в § 2.1.

§ 4.1—4.7 потребуют от читателя как знания курса анализа, так и знакомства с некоторыми элементами теории функций комплексной переменной, в частности с теорией вычетов.

Значительно большим запасом знаний нужно обладать для чтения § 4.8, где излагаются некоторые результаты исследований сходимости интерполяционных процессов. Здесь используются теоремы о сходимости последовательности линейных операторов и теорема Чебышева об альтернансе для многочленов наилучшего приближения.

Нужные сведения можно найти в § 2 добавления I и в добавлении III.

Кроме того, при изучении условий сходимости интерполирования аналитических функций необходимо иметь понятие об интегралах Стильеса и о простейших свойствах логарифмического потенциала.

Для чтения § 5.1—5.8 и 5.10 достаточно иметь сведения в объеме курса анализа.

§ 5.9 требует знания чисел и многочленов Бернулли. Все нужные сведения о них можно найти в добавлении II.

Для чтения § 5.10, где содержатся теоремы о сходимости квадратурных процессов, необходимо знание условий сходимости последовательности линейных операторов, которые можно найти в § 2 добавления I и, кроме того, нужно иметь представление о простейших свойствах интеграла Лебега.

Наконец, для понимания признаков устойчивости правил вычисления неопределенных интегралов (§ 5.11—5.13) нужно иметь простейшие све-

дения о линейных разностных уравнениях. Их можно найти в добавлении IV.

В книге § 2.1, 2.2 и 3.1—3.4 написал *В. В. Бобков*, § 2.3—2.6 и 3.5—3.7 — *П. И. Монастырный*, весь остальной текст — *В. И. Крылов*.

Авторы приносят глубокую благодарность научному редактору книги докт. физ.-мат. наук проф. И. П. Мысовских, рецензентам акад. А. Н. Тихонову, докт. физ.-мат. наук проф. М. К. Гавурину и канд. физ.-мат. наук доц. И. К. Даугавету за ценные советы и замечания, способствовавшие улучшению книги.

Авторы

Глава 1

РЕШЕНИЕ ЧИСЛЕННЫХ УРАВНЕНИЙ

§ 1.1. О СОДЕРЖАНИИ ЗАДАЧИ РЕШЕНИЯ УРАВНЕНИЙ

Задача решения уравнений в общем виде имеет указываемый ниже смысл. Пусть даны множество X элементов x и множество Y , элементы которого обозначим y . Природа элементов каждого из множеств может быть любой: это могут быть числа, совокупности чисел, функции, точки, линии и т. д. Мы не налагаем также никаких ограничений на свойства множеств X и Y и считаем их произвольными. Допустим, кроме того, что на множестве X определен оператор $y = A(x)$, который ставит в соответствие каждому элементу x из X некоторый элемент y из Y . Часто говорят, что оператор A отображает множество X в множество Y . Элемент x называют оригиналом, а $y = A(x)$ — изображением x .

Возьмем какой-либо элемент y_0 , принадлежащий Y , и поставим себе целью найти такие элементы $x \in X$, для которых y_0 является изображением. Такая задача равносильна решению операторного уравнения

$$f(x) = y_0. \quad (1.1.1)$$

Для него могут быть поставлены следующие первые проблемы.

1. Имеет ли уравнение (1.1.1) решение, т. е. существует ли такой элемент x , изображением которого будет y_0 ?

2. Если уравнение имеет решения, то при выполнении каких условий решение будет единственным? Если же решений несколько, то каким будет множество всех решений?

3. Нужно указать правило, следуя которому, можно было бы найти, в зависимости от поставленной цели и условий, точно или приближенно все решения (1.1.1), или какое-либо одно решение, заранее указанное, или любое из числа существующих.

Две первые проблемы принадлежат общей теории уравнений. В теории вычислительных методов изучается преимущественно третья из этих проблем — проблема эффективного нахождения решения уравнения.*)

*) Во многих вопросах не существует строгого разграничения между методами общей и вычислительной теории уравнений. Некоторые вычислительные методы часто применяются для доказательства существования решения уравнений. Пример этого дает метод Эйлера в теории обыкновенных дифференциальных уравнений. Его применение приводит к весьма общим теоремам о разрешимости задачи с начальными условиями. Но вместе с тем он нередко применяется и для вычислений. Другим примером может служить метод конечных разностей в уравнениях с частными производными, являющийся одним из основных методов решения задач прикладного характера, приводящих к таким уравнениям. Он одновременно часто применяется при исследовании вопросов разрешимости.

Она будет рассматриваться в частных постановках, о которых пойдет речь ниже.

Особое значение для нас будут иметь уравнения, в которых x и y будут численными величинами, X и Y — множествами их значений, а оператором, отображающим X в Y , будет некоторая функция. Уравнение (1.1.1) в этом случае можно записать в виде

$$f(x) = 0. \quad (1.1.2)$$

Мы будем рассматривать лишь численные методы решения таких уравнений и оставим в стороне многие методы, основанные на геометрическом, механическом, электрическом и других моделированиях уравнения (1.1.2).

В теории численных методов стремятся построить вычислительный процесс, при помощи которого можно найти решение (1.1.2) с наперед указанной точностью. Особенно большое значение имеют сходящиеся процессы, позволяющие решать уравнение с любой, сколь угодно малой погрешностью.

Изучение численных уравнений не является единственной основной задачей вычислительной теории уравнений. Не меньшее место в ней занимает проблема приближенного приведения операторных нечисленных уравнений к численным, что возможно сделать в большом числе случаев.

Поясним эту мысль простым примером, который, может быть, скажет лицам, приступающим к изучению вычислительных методов, больше, чем общие соображения. Пусть на отрезке *) $[a, b]$ рассматривается следующая граничная задача для дифференциального уравнения второго порядка:

$$\begin{aligned} L(x) &= x'' + p(t)x' + q(t)x = f(t), \\ a \leq t \leq b, \quad x(a) &= 0, \quad x(b) = 0, \end{aligned} \quad (1.1.3)$$

где $p(t)$, $q(t)$ и $f(t)$ предполагаются непрерывными на $[a, b]$.

За множество X , на котором определен дифференциальный оператор $L(x)$, может быть принято множество функций $x(t)$, заданных на отрезке $a \leq t \leq b$, дважды непрерывно дифференцируемых там и удовлетворяющих поставленным граничным условиям.

Предположим, что граничная задача имеет единственное решение **) и перед нами поставлен вопрос о численном нахождении значений функ-

*) Отрезок обозначается двумя буквами, которыми названы его концы, поставленными в скобках, со следующим правилом их употребления. Замкнутый конец отрезка отмечается квадратной скобкой, открытый конец — круглой, и произвольный конец — угловой скобкой. Например, если отрезок с концами a и b замкнут слева и имеет произвольный правый конец, то он обозначается знаком $[a, b >]$.

**) Так заведомо будет, если коэффициент $q(t)$ имеет отрицательные значения всюду на $[a, b]$.

ции $x(t)$. Вычислить же значения x можно только в конечном числе точек. Для решения поставленной задачи этого оказывается достаточно, так как если мы будем знать значения $x(t)$ с большой точностью на густой сетке точек отрезка $[a, b]$, то $x(t)$ можно вычислить в любой точке этого отрезка с хорошей точностью. Для таких вычислений будут даны правила, например в главе об интерполировании функций.

Для простоты рассмотрим на $[a, b]$ сетку равноотстоящих точек $t_k = a + hk$ ($h = \frac{b-a}{n}$, $k=0, 1, \dots, n$) и будем искать значения x в узлах этой сетки: $x(t_k) = x_k$. Положим в уравнении $t = t_k$ и заменим значения производных $x'(t_k)$ и $x''(t_k)$ следующими их приближенными выражениями:

$$x'(t_k) \approx \frac{x_{k+1} - x_{k-1}}{2h}, \quad x''(t_k) \approx \frac{x_{k+1} - 2x_k + x_{k-1}}{h^2}.$$

Это дает возможность дифференциальную граничную задачу (1.1.3) заменить линейной системой уравнений

$$\begin{aligned} x_{k+1} - 2x_k + x_{k-1} + \frac{h}{2} p_k (x_{k+1} - x_{k-1}) + h^2 q_k x_k &= h^2 f_k \\ (k=1, 2, \dots, n-1), \quad p_k &= p(t_k), \quad q_k = q(t_k), \quad f_k = f(t_k), \\ x_0 &= 0, \quad x_n = 0. \end{aligned} \quad (1.1.4)$$

Так как погрешность замены производных y' и y'' их выражениями через значения y_{k-1} , y_k , y_{k+1} имеет малую величину,^{*)} то можно ожидать, что решение алгебраической системы (1.1.4) будет близким к значениям $x(t_k)$ точного решения граничной задачи (1.1.3) в узлах сетки, и мы вправе принять решение системы (1.1.4) за приближенное представление решения дифференциальной граничной задачи.

Мы привели сейчас простой пример сведения операторного уравнения, где неизвестным элементом была функция $x(t)$, к системе численных уравнений. В других, более сложных задачах такое сведение часто представляет трудную проблему и хорошее решение ее может потребовать глубокого проникновения во внутреннее содержание вопроса и большой изобретательности.

Укажем на одно требование, которое обычно предъявляется к выбору метода сведения во всякой задаче и удовлетворить которое иногда бывает затруднительно. Оно связано с тем, что сложность решения численной системы быстро возрастает с увеличением количества уравнений. Поэтому при выборе способов сведения следует заботиться о том, чтобы

^{*)} Погрешность будет величиной порядка h^2 .

полученная система могла дать нужную точность по возможности при небольшом количестве численных уравнений.

Выше были указаны две основные задачи теории вычислительных методов решения уравнений: алгоритмическая теория численных уравнений и приведение нечисленных операторных уравнений к системам численных уравнений. В настоящей главе мы будем рассматривать почти исключительно численные уравнения. Задача же приведения нечисленных уравнений к численным будет изучаться в дальнейших главах только для отдельных видов уравнений — интегральных и дифференциальных, обыкновенных и с частными производными.

Общих же операторных уравнений мы кратко коснемся в настоящей главе только с целью показать, что некоторые методы решения уравнений, например итерации и Ньютона, которые мы будем изучать для численных уравнений, имеют более общее значение и могут с успехом применяться к весьма широким классам операторных уравнений.

§ 1.2. МЕТОД ИТЕРАЦИИ. СЛУЧАЙ ОДНОГО ЧИСЛЕННОГО УРАВНЕНИЯ

Общая теория метода итерации или метода повторных подстановок будет кратко изложена в следующих параграфах, сейчас же мы ознакомимся с основами теории метода на примере простейшего случая одного численного уравнения.

Выполнение итераций требует приведения уравнения к некоторой канонической форме. Допустим, что такое приведение выполнено и уравнение нам дано в виде

$$x = \varphi(x). \quad (1.2.1)$$

При этом должно быть указано множество значений, которые может принимать переменная x . Его мы обозначим X . В прикладных задачах чаще всего X будет либо вся числовая ось, либо некоторый отрезок ее. Функция $y = \varphi(x)$ каждому значению x ставит в соответствие некоторое число y .

Множество всех y , которое мы обозначим Y , образует область значений функции φ . Зависимость $y = \varphi(x)$ можно рассматривать как оператор, преобразующий X в Y .

Уравнение (1.2.1) означает, что в множестве X нужно найти такие значения x , которые переходят в себя при преобразовании оператором φ и являются, следовательно, неподвижными точками преобразования.

Задаче решения (1.2.1) легко придать геометрический смысл. В плоскости введем декартову систему координат (x, y) и построим в ней график левой части уравнения: $y = x$, являющийся биссектрисой координатного угла. Затем построим график правой части: $y = \varphi(x)$. Им будет, вообще говоря, некоторая линия плоскости, которую мы назовем l .

Решением уравнения (1.2.1) будет абсцисса точки пересечения линии l с биссектрисой $y=x$. Таких точек и соответствующих им решений может быть несколько.

Предположим, что каким-либо способом нами задано исходное приближение x_0 к решению уравнения. Все дальнейшие приближения строятся по единообразному правилу: за следующее приближение x_{n+1} принимается результат подстановки предыдущего приближения x_n в правую часть $\varphi(x)$ уравнения (1.2.1):

$$x_{n+1} = \varphi(x_n) \quad (n=0, 1, \dots). \quad (1.2.2)$$

Такое правило построения мы будем называть *простой одношаговой итерацией*.

Геометрическая картина построения приближений x_n указана на рис. 1.2.1. По исходному приближению x_0 на линии l находим точку $M_0[x_0, \varphi(x_0)]$. Через нее проводим прямую, параллельную оси x , и берем точку пересечения этой прямой с биссектрисой $y=x$. Абсциссу пересечения принимаем за x_1 и находим на l точку $M_1[x_1, \varphi(x_1)]$ и т. д.

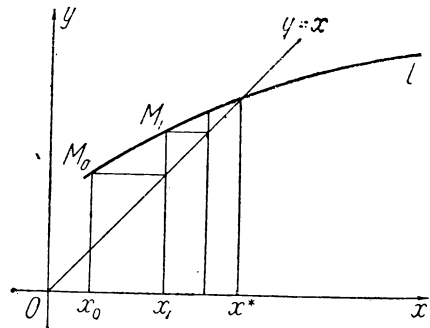


Рис. 1.2.1

x_{n+1} может быть построено, когда x_n принадлежит множеству X , на котором определена функция φ . Поэтому бесконечный итерационный процесс (1.2.2) возможен лишь в том случае, когда все x_n ($n=0, 1, \dots$) принадлежат множеству X . Это наверно будет так, если x_0 взято из X и множество Y значений $\varphi(x)$ содержится в X .*) Допустим, что итерационная последовательность $x_0, x_1, \dots, x_n, \dots$ может быть построена. В связи с ее изучением, так же, как в связи с изучением любой другой последовательности приближений к разыскиваемой величине, возникают следующие первые вопросы.

1. При каких условиях можно гарантировать возможность построения последовательности x_n ?

2. Каковы условия сходимости этой последовательности?

3. Если последовательность сходится: $\lim x_n = x^*$, то будет ли x^* решением уравнения?

4. Какова скорость сходимости, или, что равносильно, как может быть оценена разность $|x^* - x_n|$?

Несколькими страницами ниже будет доказана теорема, дающая ответ на сформулированные вопросы. Сейчас же мы остановимся на выяснении картины поведения приближений x_n вблизи решения x^* .

*) Иначе говоря, оператор φ преобразует множество X в себя.

Это позволит сделать наглядными некоторые стороны теоремы о сходимости.

Предположим, что x_n и x_{n+1} близки к решению x^* и разности $x^* - x_n = \varepsilon_n$, $x^* - x_{n+1} = \varepsilon_{n+1}$ являются малыми величинами. Допустим, кроме того, что φ имеет непрерывную производную в окрестности x^* , где лежат x_n и x_{n+1} .

Зависимость между ε_n и ε_{n+1} получится, если внести в правило итерации (1.2.2) вместо x_n и x_{n+1} их выражения через ε_n и ε_{n+1} :

$$x^* - \varepsilon_{n+1} = \varphi(x^* - \varepsilon_n) = \varphi(x^*) - \varepsilon_n \varphi'(x^*) + o(\varepsilon_n)$$

или, ввиду $x^* = \varphi(x^*)$,

$$\varepsilon_{n+1} = \varphi'(x^*) \varepsilon_n + o(\varepsilon_n). \quad (1.2.3)$$

Рассмотрим сначала простейший случай, когда $\varphi'(x^*) \neq 0$. Если ε_n — достаточно малая величина, зависимость ε_{n+1} от ε_n будет определяться приближенным равенством $\varepsilon_{n+1} \approx \varphi'(x^*) \varepsilon_n$. При $|\varphi'(x^*)| > 1$ погрешность ε_{n+1} по абсолютному значению будет больше $|\varepsilon_n|$ и приближение x_{n+1} будет отстоять от x^* дальше, чем x_n . В этом случае решение x^* будет точкой отталкивания для итерационной последовательности и поэтому здесь трудно ожидать сходимости x_n к x^* .

Если же $|\varphi'(x^*)| < 1$, то $|\varepsilon_{n+1}|$ будет меньше, чем $|\varepsilon_n|$. Поэтому можно ожидать, что если x_0 взято достаточно близко к x^* , то итерационная последовательность будет сходиться к решению. Сходимость будет происходить приблизительно по закону геометрической прогрессии со знаменателем $\varphi'(x^*)$.

Заметим еще, что в случае $\varphi'(x^*) > 0$ ε_n и ε_{n+1} , если они достаточно малы, будут иметь одинаковые знаки и x_n , начиная с некоторого места, будет стремиться монотонно к x^* . При $\varphi'(x^*) < 0$ знаки ε_n и ε_{n+1} , начиная с некоторого n , будут противоположными и сходимость x_n к x^* будет связана с колебаниями x_n около x^* . Это обстоятельство облегчает суждение о точности приближений, так как точное решение будет лежать между приближениями соседних номеров x_n и x_{n+1} .

Предположим теперь, что $\varphi'(x^*) = 0$. Тогда $\varepsilon_{n+1} = o(\varepsilon_n)$ и погрешности следующих приближений будут малыми величинами высшего порядка малости сравнительно с ε_n . В этом исключительном случае можно ожидать, что при x_0 , достаточно близком к x^* , итерационная последовательность окажется обязательно сходящейся к x^* , причем сходимость $x_n \rightarrow x^*$ будет весьма быстрой — быстрее сходимости геометрической прогрессии со сколь угодно малым знаменателем. Этим пользуются для улучшения сходимости итерационной последовательности при помощи предварительного преобразования заданного уравнения к новому $x = \psi(x)$, которое имеет то же самое решение x^* и для которого $\varphi'(x^*) = 0$. О двух видах таких преобразований мы будем говорить в следующем параграфе.

Возвратимся к соотношению между погрешностями: $\varepsilon_{n+1} = o(\varepsilon_n)$. Порядок малости ε_{n+1} зависит от кратности, с которой φ' обращается в нуль в точке x^* . Допустим, что φ имеет в окрестности x^* непрерывную производную порядка m и при этом

$$\varphi'(x^*) = \dots = \varphi^{(m-1)}(x^*) = 0, \quad \varphi^{(m)}(x^*) \neq 0.$$

Разложение $\varphi(x_n) = \varphi(x^* - \varepsilon_n)$ по степеням ε_n будет иметь вид

$$\varphi(x^* - \varepsilon_n) = \varphi(x^*) + \frac{(-1)^m}{m!} \varepsilon_n^m \varphi^{(m)}(x^*) + o(\varepsilon_n^m).$$

Подстановка его в (1.2.2) вместо $\varphi(x_n)$ приведет к следующему соотношению между погрешностями:

$$\varepsilon_{n+1} = \frac{(-1)^{m-1}}{m!} \varphi^{(m)}(x^*) \varepsilon_n^m + o(\varepsilon_n^m). \quad (1.2.4)$$

Как видно отсюда, ε_{n+1} будет малой величиной порядка m сравнительно с ε_n .*) Когда ε_n является настолько малой величиной, что можно пренебречь в (1.2.4) вторым членом правой части сравнительно с первым, (1.2.4) перейдет в приближенное равенство

$$\frac{\varepsilon_{n+1}}{\varepsilon_n^m} \approx \frac{(-1)^{m-1}}{m!} \varphi^{(m)}(x^*)$$

и, следовательно, отношение $\frac{\varepsilon_{n+1}}{\varepsilon_n^m}$ будет величиной, почти не зависящей от номера n . Отсюда сразу же вытекает, что

$$\frac{\varepsilon_{n+1}}{\varepsilon_n} \approx \left(\frac{\varepsilon_n}{\varepsilon_{n-1}} \right)^m.$$

Последнее же говорит о том, что при сделанных нами предположениях о производных функции φ , когда итерационные приближения x_n становятся близкими к решению x^* , отношение погрешностей $\frac{\varepsilon_n}{\varepsilon_{n-1}}$ после одной итерации приближенно возводится в степень m .**)

*) Если в некотором процессе последовательных приближений погрешности ε_n и ε_{n+1} связаны между собой равенством вида $|\varepsilon_{n+1}| = C|\varepsilon_n|^k + o(\varepsilon_n^k)$, $C \neq 0$, то говорят, что процесс имеет сходимость степени k .

Если $\varphi'(x^*) \neq 0$, то имеет место равенство (1.2.3) и при $|\varphi'(x^*)| < 1$ итерационный процесс имеет сходимость первой степени. В случае равенства (1.2.4) процесс приближений имеет сходимость степени m .

**) Например, если ε_n в 10 раз меньше ε_{n-1} , то ε_{n+1} будет приблизительно в 10^m раз меньше ε_n .

Предыдущие соображения очень наглядны, но неточны во многих отношениях. Например, в них не указано, сколь близким к точному решению x^* должно быть взято исходное приближение x_0 , чтобы можно было гарантировать сходимость x_n к x^* при выполнении $|\varphi(x^*)| < 1$; нет точных оценок быстроты сходимости и т. д. Ниже будет доказана одна из наиболее простых теорем, лишенная указанных недостатков. Она полезна не только для суждения о строгих достаточных условиях сходимости процесса итерации, но в ней даны также условия существования решения уравнения (1.2.1) на некотором отрезке около x_0 .

Теорема 1. Пусть выполняются условия:

1) функция $\varphi(x)$ определена на отрезке

$$|x - x_0| \leq \delta; \quad (\alpha)$$

2) непрерывна там и удовлетворяет условию Липшица с коэффициентом, меньшим единицы: *)

$$|\varphi(x) - \varphi(x')| \leq q |x - x'| \quad (0 \leq q < 1); \quad (1.2.5)$$

3) для начального значения x_0 верно неравенство

$$|x_0 - \varphi(x_0)| \leq m; \quad (1.2.6)$$

4) для чисел δ , q и m выполнено требование

$$\frac{m}{1-q} \leq \delta. \quad (1.2.7)$$

Тогда:

1) уравнение (1.2.1) на отрезке (α) имеет решение;

2) итерационная последовательность приближений может быть построена, принадлежит отрезку (α) и является сходящейся: $\lim x_n = x^*$, при этом предел x^* последовательности есть решение уравнения (1.2.1);

3) для x_n выполняется неравенство

$$|x^* - x_n| \leq \frac{m}{1-q} q^n.$$

Доказательство. Наглядный смысл теоремы весьма прост. Функция $y = \varphi(x)$ преобразует отрезок $[x_0 - \delta, x_0 + \delta]$ числовой оси в некоторый отрезок той же оси. В условии (2) x и x' есть две любые точки

*) Условие (2) часто заменяют другим: функция $\varphi(x)$ всюду на отрезке (α) имеет производную и для нее справедливо неравенство $|\varphi'(x)| \leq q < 1$. Если выполняется это последнее условие, то будет выполнено и неравенство (1.2.5), что следует из теоремы Лагранжа о приращении функции.

отрезка (α) , $|x-x'|$ — расстояние между ними. $|\varphi(x)-\varphi(x')|$ есть расстояние между точками, в которые перейдут x, x' после преобразования оператором φ . Отношение $\frac{|\varphi(x)-\varphi(x')|}{|x-x'|}$ имеет смысл «коэффициента

увеличения» расстояния при преобразовании. Неравенство Липшица (1.2.5) означает, что для любых пар точек x, x' из отрезка (α) «коэффициент увеличения» расстояния ограничен числом q . Условие $q < 1$ говорит о том, что на самом деле при преобразовании происходит уменьшение расстояний между точками по меньшей мере в q раз и отображение $y=\varphi(x)$ будет «сжатием».

Число m , входящее в неравенство (1.2.6), связано с удаленностью исходного приближения x_0 от решения уравнения x^* . Если случайно окажется $x_0=x^*$, то будет $x_0=\varphi(x_0)$ и m можно положить равным нулю. Когда же x_0 будет близким к x^* , то за m может быть взято малое число.

Условие (4) налагает ограничения на значения δ, q, m и говорит о том, что если сжатие при отображении достаточно сильное и q не близко к единице, а это верно для достаточно больших δ , т. е. в достаточно широкой окрестности около x_0 , а x_0 взято близким к решению x^* , то верны утверждения теоремы. Мерой же всех ограничений, налагаемых на δ, q, m , является неравенство (1.2.7).

Покажем сначала, что приближение x_n любого номера может быть построено, принадлежит отрезку (α) и для приближений соседних номеров выполняется неравенство

$$|x_{n+1}-x_n| \leq m q^n. \quad (1.2.8)$$

Для x_0 и x_1 это просто проверяется, так как x_0 принадлежит отрезку (α) и $x_1=\varphi(x_0)$ имеет смысл по условию (1). Далее, $|x_0-x_1| = |x_0-\varphi(x_0)| \leq m$ и неравенство (1.2.8) для x_0 и x_1 верно. Наконец, так как $m \leq \frac{m}{1-q} \leq \delta$, то $|x_0-x_1| \leq \delta$ и x_1 принадлежит (α) .

Предположим теперь, что x_0, x_1, \dots, x_n могут быть построены, принадлежат (α) и

$$|x_{k+1}-x_k| \leq m q^k \quad (k=0, 1, \dots, n-1).$$

По индуктивному предположению x_n принадлежит (α) и, так как $\varphi(x)$ определена в (α) , $x_{n+1}=\varphi(x_n)$ может быть построено. Ввиду условия (2) теоремы,

$$|x_{n+1}-x_n| = |\varphi(x_n)-\varphi(x_{n-1})| \leq q|x_n-x_{n-1}|.$$

Но для x_{n-1}, x_n верно $|x_n-x_{n-1}| \leq m q^{n-1}$ и, следовательно, $|x_{n+1}-x_n| \leq m q^n$, что доказывает для x_n и x_{n+1} неравенство (1.2.8). Наконец,

$$|x_{n+1}-x_0| = |(x_{n+1}-x_n) + (x_n-x_{n-1}) + \dots + (x_1-x_0)| \leq \\ \leq mq^n + mq^{n-1} + \dots + m = \frac{m-mq^{n+1}}{1-q} \leq \frac{m}{1-q} \leq \delta$$

и x_{n+1} принадлежит отрезку (α) . Этим закончена индукция.

Для доказательства сходимости достаточно убедиться в том, что для последовательности x_n ($n=0, 1, \dots$) выполняется признак Больцано — Коши

$$|x_{n+p}-x_n| = |(x_{n+p}-x_{n+p-1}) + (x_{n+p-1}-x_{n+p-2}) + \dots + (x_{n+1}-x_n)| \leq \\ \leq mq^{n+p-1} + mq^{n+p-2} + \dots + mq^n = \frac{mq^n - mq^{n+p}}{1-q} \leq \frac{m}{1-q} q^n, \\ |x_{n+p}-x_n| \leq \frac{m}{1-q} q^n. \quad (1.2.9)$$

Так как $q < 1$ и правая часть неравенства не зависит от p , отсюда следует выполнение признака Больцано — Коши. Стало быть, существует

$$\lim_{n \rightarrow \infty} x_n = x^*.$$

Кроме того, все x_n принадлежат замкнутому отрезку (α) , поэтому и x^* принадлежит (α) .

Вернемся к (1.2.9) и перейдем в нем к пределу при $p \rightarrow \infty$. При этом $x_{n+p} \rightarrow x^*$ и в пределе получится неравенство

$$|x^* - x_n| \leq \frac{m}{1-q} q^n,$$

доказывающее справедливость утверждения теоремы о скорости сходимости.

Осталось убедиться в том, что x^* есть решение (1.2.1).

Рассмотрим правило итерации (1.2.2) и допустим, что n неограниченно возрастает. Тогда будет $x_{n+1} \rightarrow x^*$, $x_n \rightarrow x^*$. Так как x^* принадлежит (α) и $\varphi(x)$ непрерывна в точке x^* , $\varphi(x_n)$ будет стремиться к $\varphi(x^*)$. В пределе получится $x^* = \varphi(x^*)$ и x^* действительно удовлетворяет (1.2.1).

Сделаем еще добавление о единственности решения.

Теорема 2. Уравнение $x = \varphi(x)$ на всяком множестве точек, на котором $\varphi(x)$ выполняет неравенство

$$|\varphi(x) - \varphi(y)| < |x - y| \quad (x \neq y),$$

может иметь не больше одного решения.

Доказательство. Пусть x и y принадлежат такому множеству и удовлетворяют уравнению $x = \varphi(x)$ и $y = \varphi(y)$. Оценим разность $x - y$, полагая $x \neq y$:

$$|x - y| = |\varphi'(x) - \varphi(y)| < |x - y|.$$

Последнее неравенство при $x \neq y$ выполняться не может, и мы должны считать $x = y$. Двух различных решений быть не может.

Теорема 1 указывает условия, достаточные для существования решения x^* в окрестности начального приближения x_0 и сходимости итерационной последовательности к x^* не медленнее, чем показывает оценка.

$$|x^* - x_n| \leq \frac{m}{1 - q} q^n.$$

Мы обращали внимание на то, что в исключительных случаях, не предусмотренных в теореме, когда $\varphi'(x^*) = 0$, сходимость может быть значительно более быстрой. Полезно дополнить теорему 1 замечанием об оценке скорости сходимости $x_n \rightarrow x^*$ в этих исключительных случаях.

Предположим, что на некотором отрезке $|x - x^*| \leq \delta$ около x^* , функция $\varphi(x)$ имеет непрерывную производную порядка m , при этом

$$\varphi'(x^*) = \dots = \varphi^{(m-1)}(x^*) = 0 \quad \text{и} \quad |\varphi^{(m)}(x)| \leq M_m,$$

когда $|x - x^*| \leq \delta$. Допустим также, что, начиная с некоторого номера N , итерационные приближения x_n ($n \geq N$) все лежат на отрезке $|x - x^*| \leq \delta$. Так как нумерацию приближений мы можем начать с любого места последовательности, допустимо для упрощения записи считать $N = 0$ и все приближения x_n — принадлежащими указанному отрезку.

Если из равенства $x^* = \varphi(x^*)$ вычесть почленно рекурсионное равенство $x_{n+1} = \varphi(x_n)$, мы получим следующую связь между погрешностями $\varepsilon_n = x^* - x_n$ двух соседних номеров:

$$\begin{aligned} \varepsilon_{n+1} &= \varphi(x^*) - \varphi(x_n) = -[\varphi(x^* - \varepsilon_n) - \varphi(x^*)] = \\ &= -\frac{(-1)^m}{m!} \varphi^{(m)}(x^* - \Theta \varepsilon_n) \varepsilon_n^m \quad (0 < \Theta < 1). \end{aligned}$$

Отсюда получается нужная нам оценка

$$|\varepsilon_{n+1}| \leq \frac{M_m}{m!} |\varepsilon_n|^m. \quad (1.2.10)$$

Если это неравенство применить n раз, начиная с ε_n , найдем следующую оценку погрешности:

$$|\varepsilon_n| \leq \alpha^{-\frac{1}{m-1}} (\alpha^{-\frac{1}{m-1}} \varepsilon_0)^{m^n} \quad \left(\alpha = \frac{M_m}{m!} \right). \quad (1.2.11)$$

§ 1.3. О ЗАДАЧЕ УЛУЧШЕНИЯ МЕТОДА ИТЕРАЦИИ. НЕКОТОРЫЕ ВИДОИЗМЕНЕНИЯ ИТЕРАЦИОННОГО ПРОЦЕССА

В предыдущем параграфе, где рассматривался метод простой одношаговой итерации, мы обратили внимание на то, что если итерационная последовательность x_n ($n=0, 1, \dots$) лежит вблизи решения x^* , то погрешности $x^* - x_n = \varepsilon_n$ изменяются, вообще говоря, приблизительно по закону геометрической прогрессии со знаменателем $\varphi'(x^*)$:

$$\varepsilon_{n+1} \approx \varphi'(x^*) \varepsilon_n.$$

Последовательность x_n будет сходиться к решению x^* , если в окрестности x^* производная φ' будет по абсолютной величине меньше единицы и если исходное приближение x_0 взято достаточно близко к x^* . При этом сходимость будет тем быстрее, чем меньшее значение имеет $|\varphi'(x^*)|$. Если же значение $|\varphi'(x^*)|$ близко к единице, сходимость $x_n \rightarrow x^*$ может быть весьма медленной и может потребоваться много итераций, чтобы достигнуть нужной точности в вычислении x^* .

Метод итерации, как и всякий другой процесс приближений, можно пытаться усовершенствовать в двух направлениях: во-первых, улучшить скорость сходимости и, во-вторых, расширить область применимости, т. е. сделать процесс сходящимся при менее ограничительных условиях, чем указывалось, например, в теореме 1 § 1.2.

Достигнуть этого можно либо заменой итеративного процесса (1.2.2) другим более сложным процессом, с более быстрой сходимостью (два видоизменения такого рода будут рассмотрены в настоящем параграфе), либо предварительно преобразовав заданное уравнение $x = \varphi(x)$ к такому виду, для которого простой одношаговый процесс сходится быстрее, чем для заданного уравнения. По этому поводу мы заметим прежде всего, что если нам дано уравнение $f(x) = 0$, то привести его к каноническому для метода итерации виду $x = \varphi(x)$ можно обычно многими способами и среди возможных способов вычислитель должен избрать тот, в котором производная $\varphi'(x)$ вблизи разыскиваемого корня имеет возможно малое абсолютное значение. При приведении используются индивидуальные свойства каждого уравнения и никаких общих правил здесь, по-видимому, дать невозможно. Успех зависит почти исключительно от опыта и искусства вычислителя. Мы оставим этот вопрос в стороне и обратим внимание на другой возможный путь преобразований.^{*)}

Напомним, что если $\varphi(x)$ имеет в окрестности решения x^* непрерывные производные порядка m и $\varphi'(x^*) = \dots = \varphi^{(m-1)}(x^*) = 0$, а $\varphi^{(m)}(x^*) \neq 0$, то погрешности ε_{n+1} и ε_n связаны равенством вида (1.2.4) и, когда x_{n+1} и x_n лежат достаточно близко к x^* , порядок малости ε_{n+1} будет в m раз выше порядка ε_n и можно ожидать весьма быстрой сходимости $x_n \rightarrow x^*$. Это приводит к мысли заменить заданное уравнение $x = \varphi(x)$ новым

^{*)} Некоторые сведения для линейных систем по этому вопросу приведены в гл. 2.

уравнением $x = \psi(x)$, которое имеет то же решение x^* , что и заданное, но для которого

$$\psi'(x^*) = \dots = \psi^{(m-1)}(x^*) = 0, \quad \psi^{(m)}(x^*) \neq 0.$$

С некоторыми способами составления уравнения $x = \psi(x)$ мы ознакомимся в следующем параграфе.

Обратимся к проблеме изменения итерационного процесса (1.2.2): $x_{n+1} = \varphi(x_n)$. Процесс одношаговый, и это является одним из его до-

x	$\varphi(x)$
\bar{x}_0	$\varphi(\bar{x}_0)$
\bar{x}_1	$\varphi(\bar{x}_1)$
\dots	\dots
\bar{x}_{n-1}	$\varphi(\bar{x}_{n-1})$
\bar{x}_n	$\varphi(\bar{x}_n)$

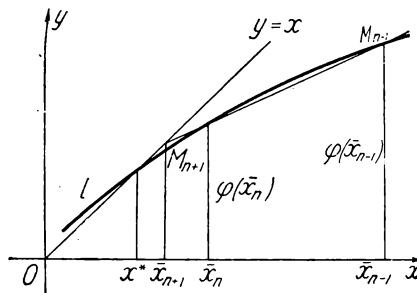


Рис. 1.3.1

стоинств, так как при его применении нужно задать только x_0 и не нужно составлять программы вспомогательного начала расчетов. Но, как во всяком одношаговом методе, мы не используем многих возможностей, которые дает вычислительный процесс. В самом деле, допустим, что выполнено n шагов итерации и составлена приводимая в тексте таблица значений приближений к решению и соответствующих значений φ . Чтобы отличить приближения, полученные по новому правилу вычислений, от приближений по правилу (1.2.2), обозначим их \bar{x}_n . Нахождению подлечит \bar{x}_{n+1} .

В правиле (1.2.2) мы используем только одно значение аналогичной таблицы, полагая $\bar{x}_{n+1} = \varphi(\bar{x}_n)$. По существу это означает, что мы заменяем функцию $\varphi(x)$ постоянной величиной $\varphi(\bar{x}_n)$, проводим через точку $M_n[\bar{x}_n, \varphi(\bar{x}_n)]$ прямую линию $y = \varphi(\bar{x}_n)$, параллельную оси x , и определяем точку пересечения ее с прямой $y = x$. Абсциссу такой точки принимаем за \bar{x}_{n+1} . Геометрическая картина процесса изображена на рис. 1.3.1.

Но для вычисления \bar{x}_{n+1} мы можем воспользоваться любыми значениями \bar{x}_k и $\varphi(\bar{x}_k)$, приведенными в таблице. Если встать на этот путь, то мы должны будем избрать тот или иной способ интерполирования \bar{x}_{n+1}

по нескольким предшествующим значениям \bar{x}_k и соответствующим им $\varphi(\bar{x}_k)$. При этом можно выиграть в скорости сходимости по сравнению с правилом (1.2.2), но мы должны будем поступиться преимуществом одношагового метода, так как новое правило вычислений будет многошаговым. Остановимся только на случае двухшагового итерационного процесса, когда выполняется линейное интерполирование $\varphi(x)$ по двум парам чисел $[\bar{x}_n, \varphi(\bar{x}_n)]$, $[\bar{x}_{n-1}, \varphi(\bar{x}_{n-1})]$. Геометрически это значит (рис. 1.3.1), что линия $l: y = \varphi(x)$ заменяется секущей прямой, проходящей через две точки $M_{n-1}[\bar{x}_{n-1}, \varphi(\bar{x}_{n-1})]$ и $M_n[\bar{x}_n, \varphi(\bar{x}_n)]$. Использование интерполирования более высокой степени мы сейчас оставим в стороне и вернемся к нему в гл. 4.

Если же мы хотим при улучшении правила (1.2.2) сохранить преимущество одношагового процесса, то увеличивать точность правила мы можем лишь за счет вычисления одного или нескольких вспомогательных значений φ . Одно из правил такого вида, принадлежащее Стеффенсену, будет рассмотрено в конце настоящего параграфа.

1. Рассмотрим метод секущих, или правило линейной интерполяции. На линии $l: y = \varphi(x)$, являющейся графиком правой части уравнения, возьмем две точки $M_n[\bar{x}_n, \varphi(\bar{x}_n)]$ и $M_{n-1}[\bar{x}_{n-1}, \varphi(\bar{x}_{n-1})]$, проведем через них секущую прямую и найдем точку пересечения ее с биссектрисой $y = x$. Абсциссу этой точки примем за \bar{x}_{n+1} . Решение системы уравнений секущей и биссектрисы:

$$\frac{(x - \bar{x}_{n-1})}{(\bar{x}_n - \bar{x}_{n-1})} = \frac{[y - \varphi(\bar{x}_{n-1})]}{[\varphi(\bar{x}_n) - \varphi(\bar{x}_{n-1})]}, \quad y = x$$

приведет к следующему правилу нахождения \bar{x}_{n+1} :

$$\bar{x}_{n+1} = \frac{\bar{x}_{n-1}\varphi(\bar{x}_n) - \bar{x}_n\varphi(\bar{x}_{n-1})}{\varphi(\bar{x}_n) - \bar{x}_n - \varphi(\bar{x}_{n-1}) + \bar{x}_{n-1}}. \quad (1.3.1)$$

Чтобы начать вычисления, необходимо указать два исходных приближения x_0, x_1 к решению.

Выясним теперь при помощи нестрогих, но наглядных соображений характер изменения погрешности $\varepsilon_k = x^* - \bar{x}_k$, когда \bar{x}_k находится вблизи решения x^* . Правило (1.3.1) дает нам связь между погрешностями приближений трех смежных номеров. Для получения ее достаточно в (1.3.1) подставить вместо \bar{x}_k значение $\bar{x}_k = x^* - \varepsilon_k$; заменить $\varphi(\bar{x}_k)$ разложением по степеням ε_k

$$\varphi(\bar{x}_k) = \varphi(x^* - \varepsilon_k) = \varphi(x^*) - \varepsilon_k \varphi'(x^*) + \frac{1}{2} \varepsilon_k^2 \varphi''(x^*) - \dots$$

$$(k = n-1, n, n+1)$$

и сохранить лишь главные члены в числителе и знаменателе дроби: *)

$$\varepsilon_{n+1} \approx - \frac{1}{2} \cdot \frac{\varphi''(x^*) \varepsilon_n \varepsilon_{n-1}}{\varphi'(x^*) - 1}. \quad (1.3.2)$$

Погрешность ε_{n+1} будет иметь, вообще говоря, тот же порядок малости, что и произведение $\varepsilon_n \varepsilon_{n-1}$, и будет поэтому малой высшего порядка сравнительно с каждой из величин ε_{n-1} и ε_n . Равенство (1.3.2) позволяет ожидать, что если $\varphi(x)$ в окрестности x^* дважды непрерывно дифференцируема, производная $\varphi'(x)$ не близка к единице в этой окрестности и, наконец, \bar{x}_0 и \bar{x}_1 взяты достаточно близко к решению x^* , то правило (1.3.1) дает последовательность приближений \bar{x}_n , сходящуюся к решению x^* со скоростью, намного превышающей скорость сходимости простого одношагового итерационного процесса (1.2.2). Сейчас мы ограничимся изложением приведенных наглядных соображений, точную же теорему о сходимости метода линейного интерполирования докажем в § 1.8, когда будем рассматривать метод секущих прямых, связанный с методом Ньютона.

2. Теперь приведем пример улучшения итерационного процесса при помощи вычисления вспомогательных значений функции. Идея метода связана с одним из способов улучшения сходимости последовательностей, изменяющихся по закону, близкому к геометрической прогрессии.

Пусть дана произвольная последовательность $s_0, s_1, s_2, \dots, s_n, \dots$. Подвергнем ее нелинейному преобразованию, носящему имя Эйткена,**) и построим новую последовательность $\sigma_1, \sigma_2, \dots$, где

$$\sigma_n = \frac{s_{n+1}s_{n-1} - s_n^2}{s_{n+1} - 2s_n + s_{n-1}}.$$

*) Более подробный анализ показал бы, что для погрешностей верно равенство

$$\varepsilon_{n+1} = - \frac{1}{2} \cdot \frac{\varphi''(\xi) \varepsilon_n \varepsilon_{n-1}}{\varphi'(\eta) - 1},$$

где ξ и η есть некоторые точки отрезка, на котором лежат x_{n-1} , x_n и x^* .

**) Задача улучшения сходимости последовательностей и рядов рассматривается во втором томе. Там же дано более подробное изложение теории преобразования Эйткена и ему родственных преобразований. Краткое описание преобразований, приведенное в настоящем параграфе, сделано только для выяснения идеи, на которой построен метод Стеффенсена.

Условием возможности преобразования является соблюдение неравенств

$$s_{n+1} - 2s_n + s_{n-1} \neq 0 \quad (n=1, 2, \dots).$$

Если

$$s_n = C + Aq^n \quad (A \neq 0, q \neq 0, 1),$$

преобразование будет возможным, так как

$$s_{n+1} - 2s_n + s_{n-1} = Aq^{n-1} (q-1)^2 \neq 0.$$

Простым подсчетом можно убедиться в том, что при любом значении n , независимо от величин q и A , s_n будет равно C . Когда $|q| < 1$, последовательность s_n будет сходиться и иметь своим пределом C . Преобразование Эйткена на любом шаге будет давать $\sigma_n = C = \lim s_n$. Если же последовательность s_n будет сходиться к C не точно по указанному выше показательному закону, а по закону, близкому к нему, то σ_n может не совпадать с C тождественно. Но можно ожидать, что σ_n будет близким к C и сходимость $\sigma_n \rightarrow C$ будет более быстрой, чем сходимость $s_n \rightarrow C$.

Напомним теперь, что простой одношаговый итерационный процесс (1.2.2), если $|\varphi'(x)| < 1$ в окрестности решения x^* , будет сходиться по закону, близкому к геометрической прогрессии со знаменателем $q = |\varphi'(x^*)|$. Для улучшения сходимости здесь естественно воспользоваться преобразованием Эйткена, изменив, однако, вычислительный процесс так, чтобы каждое вновь найденное улучшенное значение сразу же вводилось в вычисления и последующие приближения находились с учетом уже сделанного улучшения. Поясним это на одном шаге преобразования. Допустим, что мы начинаем вычисления с исходного значения x_0 . Пользуясь правилом (1.2.2), строим $x_1 = \varphi(x_0)$, $x_2 = \varphi(x_1) = \varphi[\varphi(x_0)]$ и к трем числам x_0, x_1, x_2 применяем преобразование Эйткена

$$x_1' = \frac{x_2 x_0 - x_1^2}{x_2 - 2x_1 + x_0} = \frac{x_0 \varphi[\varphi(x_0)] - \varphi^2(x_0)}{\varphi[\varphi(x_0)] - 2\varphi(x_0) + x_0}.$$

Этим мы закончим один шаг процесса Стеффенсена. Чтобы сделать второй шаг, мы выполняем, начиная с x_1' , те же вычисления, какие были нами сделаны для x_0 , и т. д.

В общем виде правило вычислений может быть сформулировано, как указано несколькими строками ниже. Приближения, найденные по этому правилу, мы обозначим x_k' ($k=0, 1, \dots$). Пусть вычисления выполнены до шага с номером n и найдены приближения x_k' ($k=0, 1, \dots, n$). Тогда приближение x'_{n+1} вычисляется по правилу

$$x'_{n+1} = \frac{x_n' \varphi[\varphi(x_n')] - \varphi^2(x_n')}{\varphi[\varphi(x_n')] - 2\varphi(x_n') + x_n'} \quad (1.3.3)$$

Правило Стеффенсена (1.3.3) является одношаговым и требует вычисления двух значений функции φ на каждый шаг.

Равенство (1.3.3) можно рассматривать как простой одношаговый процесс вида (1.2.2) для вспомогательного уравнения

$$x = \Phi(x), \quad \Phi(x) = \frac{x\varphi[\varphi(x)] - \varphi^2(x)}{\varphi[\varphi(x)] - 2\varphi(x) + x}. \quad (1.3.4)$$

Чтобы выяснить причину, в силу которой можно ожидать, что правило Стеффенсена (1.3.3) имеет лучшую сходимость, чем правило (1.2.2), мы воспользовались преобразованием Эйткена. Оно позволило сделать вполне наглядной интуитивную сторону правила (1.3.3).

Преобразование Эйткена взято нами из совсем другой области теории численных методов — из проблемы улучшения сходимости последовательностей — и на первый взгляд может показаться, что правило Стеффенсена никак не связано с линейным интерполированием функции φ и методом секущих.

Мы покажем сейчас, что такая связь все-таки существует, но линейное интерполирование, которое нужно осуществить для получения правила Стеффенсена, отличается от того, которое было применено для вывода правила секущих (1.3.1).

Пусть x_n' и $\varphi(x_n')$ нам известны и известна, следовательно, точка $M_n'[x_n', \varphi(x_n')]$ на линии $y = \varphi(x)$. Построим на этой линии вспомогательную точку $M_n''(x_n'', y_n'')$ по простому одношаговому правилу итерации, положив

$$x_n'' = \varphi(x_n'), \quad y_n'' = \varphi(x_n'') = \varphi[\varphi(x_n')].$$

Через точки M_n' и M_n'' проведем секущую прямую, уравнение которой

$$\frac{(x - x_n')}{\{\varphi(x_n') - x_n'\}} = \frac{[y - \varphi(x_n')]}{\{\varphi[\varphi(x_n')] - \varphi(x_n')\}},$$

и найдем точку ее пересечения с биссектрисой $y = x$. Абсциссу последней точки примем за следующее приближение к решению уравнения $x = \varphi(x)$ и обозначим это приближение x'_{n+1} . Чтобы получить правило его вычисления, достаточно в уравнении секущей положить $y = x = x'_{n+1}$ и из полученного после этого уравнения найти x'_{n+1} . Простые вычисления покажут, что полученное правило совпадет с правилом (1.3.3).

Выясним теперь преимущества, которые имеет в отношении скорости сходимости правило (1.3.3) сравнительно с правилом (1.2.2).

Рассмотрим сначала простейший и основной случай в поставленной задаче сравнения. Допустим, что $\varphi(x)$ имеет непрерывную производную $\varphi'(x)$ в некоторой окрестности решения $x = x^*$, при этом $\varphi'(x^*) = \alpha \neq 0$. Напомним, что в этих условиях мы можем наверное гарантировать сходимость итерационного процесса (1.2.2) к x^* только в том

случае, когда исходное приближение x_0 взято достаточно близко к x^* и $|\varphi'(x^*)| = |\alpha| < 1$. При этом сходимость будет происходить по закону (1.2.3), близкому к геометрической прогрессии.

Для итерационного процесса Стеффенсена (1.3.3) условия сходимости, как следует из приводимых ниже лемм, будут значительно более благоприятными.

Лемма 1. Если для $\varphi(x)$ вблизи x^* верно представление

$$\varphi(x) = x^* + \alpha(x - x^*) + o(x - x^*) \quad (1.3.5)$$

и $\alpha \neq 0$, $\alpha \neq 1$, то $\Phi(x) = x^* + o(x - x^*)$.

Доказательство. Не уменьшая общности, мы можем считать $x^* = 0$. Для достижения этого достаточно положить $x = x^* + z$, $\varphi(x) - x^* = \varphi(x^* + z) - x^* = \psi(z)$ и рассмотреть уравнение $z = \psi(z)$ и соответствующее ему правило Стеффенсена $z'_{n+1} = \Psi(z'_n)$:

$$\Psi(z) = \frac{z\psi[\psi(z)] - \psi^2(z)}{\psi[\psi(z)] - 2\psi(z) + z} = \Phi(z) - x^*,$$

При $x^* = 0$ равенство (1.3.5) примет форму $\varphi(x) = \alpha x + o(x)$. Тогда

$$\varphi[\varphi(x)] = \alpha[\alpha x + o(x)] + o(x) = \alpha^2 x + o(x),$$

$$\varphi[\varphi(x)] - 2\varphi(x) + x = \alpha^2 x - 2\alpha x + x + o(x) = (\alpha - 1)^2 x + o(x).$$

Отсюда, в частности, следует, что делитель в выражении (1.3.4) функции $\Phi(x)$ при x , близком к решению $x^* = 0$, и $x \neq 0$ будет отличным от нуля. Далее,

$$x\varphi[\varphi(x)] - \varphi^2(x) = \alpha^2 x^2 + o(x^2) - [\alpha x + o(x)]^2 = o(x^2),$$

$$\Phi(x) = \frac{o(x^2)}{(\alpha - 1)^2 x + o(x)} = o(x),$$

что доказывает лемму.

Если положить $\Phi(0) = 0$, то функция $\Phi(x)$ будет непрерывна в точке $x = 0$ и, кроме того, $\Phi'(0) = 0$.

Предположим теперь дополнительно, что $\varphi(x)$ имеет непрерывную производную в окрестности точки $x = x^* = 0$. Так как в выражении (1.3.4) для $\Phi(x)$ делимое и делитель будут непрерывно дифференцируемы в окрестности $x = 0$, то $\Phi(x)$ будет иметь непрерывную производную во всех точках некоторой окрестности $x = 0$, кроме, может быть, самой точки $x = 0$. Покажем, что $\Phi'(x)$ будет непрерывна и при $x = 0$. Для этого достаточно показать, что $\Phi'(x) \rightarrow 0$ ($x \rightarrow 0$).

Умножим обе части второго равенства (1.3.4) на делитель и от обеих частей возьмем производную:

$$\begin{aligned} \{ \varphi[\varphi(x)] - 2\varphi(x) + x \} \Phi'(x) + \{ \varphi'[\varphi(x)]\varphi'(x) - 2\varphi'(x) + 1 \} \Phi(x) = \\ = \varphi[\varphi(x)] + x\varphi'[\varphi(x)]\varphi'(x) - 2\varphi(x)\varphi'(x). \end{aligned} \quad (1.3.6)$$

Воспользуемся тем, что $\varphi(x) = \alpha x + o(x)$ и $\varphi'(x) = \alpha + o(1)$. Множитель при $\Phi'(x)$, как выяснилось выше, имеет при $x \rightarrow 0$ форму $(\alpha - 1)^2 x + o(x)$. Второй член левой части равенства является величиной порядка $o(x)$, так как $\Phi(x) = o(x)$, а множитель при $\Phi(x)$ ограничен. Для правой же части равенства верно представление

$$\begin{aligned} \varphi[\varphi(x)] + x\varphi'[\varphi(x)]\varphi'(x) - 2\varphi(x)\varphi'(x) = \alpha^2 x + o(x) + x[\alpha + o(1)][\alpha + o(1)] - \\ - 2[\alpha x + o(x)][\alpha + o(1)] = o(x). \end{aligned}$$

Из (1.3.6) следует

$$\{(\alpha-1)^2 x + o(x)\} \Phi'(x) = o(x),$$

что возможно только при условии $\Phi'(x) = o(1)$, так как

$$\Phi'(x) = \frac{o(1)}{(\alpha-1)^2 + o(1)} = o(1).$$

Все изложенное выше позволяет сформулировать лемму:

Лемма 2. Если функция $\varphi(x)$ имеет непрерывную производную в окрестности решения x^* уравнения $x = \varphi(x)$ и если $\varphi'(x^*) = \alpha \neq 0$ и $\alpha \neq 1$, то функция $\Phi(x)$, определенная равенством (1.3.4) и дополнительным условием $\Phi(x^*) = x^*$, будет непрерывно дифференцируемой в некоторой окрестности решения x^* , при этом $\Phi'(x^*) = 0$.

Лемма позволяет утверждать, что если x_0' взято достаточно близким к решению x^* , то последовательность x_m' , построенная по правилу Стеффенсена (1.3.3), сходится к x^* , при этом сходимость будет настолько быстрой, что для погрешности $\varepsilon_n' = x^* - x_n'$ будет выполняться соотношение $\frac{\varepsilon_{n+1}'}{\varepsilon_n'} \rightarrow 0$ и сходимость $x_n' \rightarrow x^*$ будет более быстрой, чем сходимость геометрической прогрессии со сколь угодно малым знаменателем.

О скорости сходимости можно получить более точное представление, если известны дополнительные сведения о поведении $\varphi(x)$ вблизи решения x^* .

Лемма 3. Если для $\varphi(x)$ при x , близком к x^* , верно равенство

$$\varphi(x) = x^* + \alpha(x - x^*) + \beta(x - x^*)^m + o[(x - x^*)^m] \quad (1.3.7)$$

и $\alpha \neq 0$, $\alpha \neq 1$, $m > 1$, $\beta \neq 0$, то

$$\Phi(x) = \frac{\alpha\beta(\alpha^m - 1)}{(\alpha - 1)^2} x^m + o(x^m). \quad (1.3.8)$$

Доказательство. Вновь будем считать $x^* = 0$. Приводимые ниже вычисления не требуют пояснений.

$$\varphi(x) = \alpha x + \beta x^m + o(x^m),$$

$$\begin{aligned} \varphi[\varphi(x)] &= \alpha[\alpha x + \beta x^m + o(x^m)] + \beta[\alpha x + \beta x^m + o(x^m)]^m + o(x^m) = \\ &= \alpha^2 x + (\alpha\beta + \beta\alpha^m) x^m + o(x^m); \end{aligned}$$

$$\begin{aligned} \varphi[\varphi(x)] - 2\varphi(x) + x &= \alpha^2 x + (\alpha\beta + \beta\alpha^m) x^m + o(x^m) - 2[\alpha x + \beta x^m + o(x^m)] + x = \\ &= (\alpha - 1)^2 x + (\alpha\beta + \beta\alpha^m - 2\beta) x^m + o(x^m); \end{aligned}$$

$$\begin{aligned} x\varphi[\varphi(x)] - \varphi^2(x) &= \alpha^2 x^2 + (\alpha\beta + \beta\alpha^m) x^{m+1} + o(x^{m+1}) - [\alpha x + \beta x^m + o(x^m)]^2 = \\ &= (\beta\alpha^m - \alpha\beta) x^{m+1} + o(x^{m+1}), \end{aligned}$$

$$\Phi(x) = \frac{(\beta\alpha^m - \alpha\beta) x^{m+1} + o(x^{m+1})}{(\alpha - 1)^2 x + (\alpha\beta + \beta\alpha^m - 2\beta) x^m + o(x^m)} = \frac{\beta\alpha^m - \alpha\beta}{(\alpha - 1)^2} x^m + o(x^m).$$

Лемма 4. Пусть функция $\varphi(x)$ имеет непрерывную производную порядка m в окрестности точки x^* и формула Тейлора по степеням $x - x^*$ для $\varphi(x)$ имеет вид:

$$\begin{aligned}\varphi(x) &= x^* + \varphi'(x^*)(x-x^*) + \frac{1}{m!} \varphi^{(m)}(x^*)(x-x^*)^m + o[(x-x^*)^m] = \\ &= x^* + \alpha(x-x^*) + \beta(x-x^*)^m + o[(x-x^*)^m].\end{aligned}\quad (1.3.9)$$

Тогда функция $\Phi(x)$, определенная равенством (1.3.4) и дополнительным условием $\Phi(x^*) = x^*$, будет иметь в некоторой окрестности точки x^* непрерывную производную порядка m и разложение $\Phi(x)$ по степеням $x-x^*$ имеет форму:

$$\Phi(x) = x^* + \frac{\alpha\beta(\alpha^{m-1}-1)}{(\alpha-1)^2} (x-x^*)^m + o[(x-x^*)^m]. \quad (1.3.10)$$

Доказательство. С целью упростить запись, как и выше, будем считать $x^* = 0$ и разложение (1.3.9) запишем в виде $\varphi(x) = \alpha x + \beta x^m + o(x^m)$.

Рассмотрим числитель и знаменатель в выражении (1.3.4) для $\Phi(x)$:

$$x\varphi[\varphi(x)] - \varphi^2(x) = M(x), \quad \varphi[\varphi(x)] - 2\varphi(x) + x = N(x).$$

При сделанных предположениях о $\varphi(x)$, обе функции M и N будут иметь в окрестности $x=0$ непрерывные производные до порядка m включительно. Степенные разложения для них, полученные при доказательстве третьей леммы,

$$M(x) = \beta\alpha(\alpha^{m-1}-1)x^{m+1} + o(x^{m+1})$$

и

$$N(x) = (\alpha-1)^2 x + \beta(\alpha^m + \alpha - 2)x^m + o(x^m)$$

являются по существу степенными разложениями Тейлора с остаточными членами в форме Пеано и позволяют судить о значениях производных этих функций при $x=0$.

Отметим, что при x , близких к нулю и отличных от нуля, $N(x)$ не обращается в нуль и отношение $\frac{M(x)}{N(x)} = \Phi(x)$ определяет $\Phi(x)$ в некоторой окрестности точки $x=0$, исключая саму точку $x=0$, и позволяет утверждать, что на указанном множестве точек $\Phi(x)$ является m -кратно непрерывно дифференцируемой функцией x . Что же касается точки $x=0$, то $\Phi(x)$ в ней мы определим по непрерывности. Это можно сделать, например, если в равенство $\Phi = \frac{M}{N}$ вместо M и N внести их разложения по степеням x и сократить делимое и делитель на первую степень x :

$$\Phi(x) = \frac{\alpha\beta(\alpha^{m-1}-1)x^m + o(x^m)}{(\alpha-1)^2 + \beta(\alpha^m + \alpha - 2)x^{m-1} + o(x^{m-1})} = \frac{\alpha\beta(\alpha^{m-1}-1)}{(\alpha-1)^2} x^m + o(x^m). \quad (1.3.11)$$

Из предыдущего изложения и последнего равенства следует, что $\Phi(x)$ после определения ее будет функцией, m -кратно непрерывно дифференцируемой всюду в некоторой окрестности точки $x=0$. Этим доказана лемма 4.

Она позволяет сказать, что если для $\varphi(x)$ верно разложение (1.3.9) вблизи решения x^* уравнения $x=\varphi(x)$ и если исходное приближение x_0 взято достаточно близким к x^* , то итерационный процесс Стеффенсена (1.3.3) сходится к решению x^* , при этом для погрешности $\varepsilon_n' = x^* - x_n'$ верно соотношение

$$\varepsilon'_{n+1} = (-1)^{m-1} \frac{\alpha\beta(\alpha^{m-1}-1)}{(\alpha-1)^2} (\varepsilon_n')^m + o[(\varepsilon_n')^m].$$

В частности, когда $m=2$ и разложение $\Phi(x)$ вблизи решения имеет вид

$$\Phi(x) = x^* + \Phi'(x^*)(x-x^*) + \frac{1}{2} \Phi''(x^*)(x-x^*)^2 + o[(x-x^*)^2];$$

соотношение между погрешностями двух приближений соседних номеров будет

$$\epsilon'_{n+1} = \frac{\Phi'(x^*)\Phi''(x^*)}{\Phi'(x^*)-1} \frac{(\epsilon_n')^2}{2} + o[(\epsilon_n')^2].$$

Во всем предшествующем изложении мы полагали $\alpha = \Phi'(x^*) \neq 1$. Рассмотрим теперь исключительный случай $\alpha = 1$ и покажем, что сходимость правила (1.3.3), вообще говоря, сохранится, но будет значительно медленнее, чем в предыдущих случаях, и близкой к геометрической прогрессии со знаменателем, меньшим 1.

Лемма 5. Пусть $\Phi'(x^*) = \alpha = 1$. Если $\Phi'(x)$ непрерывна в окрестности решения x^* и имеет место равенство

$$\Phi'(x) - 1 = T(x)(x-x^*)^{m-1},$$

где $m > 1$ и $T(x)$ стремится к конечному пределу $\gamma \neq 0$ при $x \rightarrow x^*$, то $\Phi(x)$ имеет производную в точке $x = x^*$ и $\Phi'(x^*) = 1 - \frac{1}{m}$. При этом считается $\Phi(x^*) = x^*$.

Доказательство. Будем по-прежнему полагать $x^* = 0$. Вычтя x из обеих частей равенства

$$\Phi(x) = \frac{x\Phi[\Phi(x)] - \Phi^2(x)}{\Phi[\Phi(x)] - 2\Phi(x) + x}$$

и обозначив $g(x) = \Phi(x) - x$, найдем

$$\Phi(x) = - \frac{g^2(x)}{g[\Phi(x)] - g(x)}.$$

По теореме о приращении функции

$$g[\Phi(x)] - g(x) = [\Phi(x) - x]g'(\xi) = g(x)g'(\xi),$$

где $\xi = x + \vartheta[\Phi(x) - x]$; $0 < \vartheta < 1$.

Так как при $x \rightarrow 0$ разность

$$\Phi(x) - x = \int_0^x T(t)t^{m-1}dt$$

будет величиной малой, порядка более высокого, чем x , x и ξ будут эквивалентны между собой: $\xi = x + o(x)$. Поэтому

$$g'(\xi) = \Phi'(\xi) - 1 = T(\xi)\xi^{m-1} = \gamma\xi^{m-1} + o(\xi^{m-1}) = \gamma x^{m-1} + o(x^{m-1}).$$

Далее,

$$\Phi(x) - x = \int_0^x [\Phi'(t) - 1]dt = \int_0^x T(t)t^{m-1}dt = \frac{\gamma}{m} x^m + o(x^m)$$

и

$$\Phi(x) = -\frac{g(x)}{g'(\xi)} = \frac{\frac{\gamma}{m}x^m + o(x^m)}{\gamma x^{m-1} + o(x^{m-1})} = -\frac{1}{m}x + o(x).$$

Наконец, ввиду $\Phi(0) = 0$,

$$\Phi'(0) = \lim_{x \rightarrow 0} \frac{1}{x} [\Phi(x) - \Phi(0)] = \lim_{x \rightarrow 0} \frac{1}{x} \left[x - \frac{1}{m}x + o(x) \right] = 1 - \frac{1}{m}.$$

Лемма 5 позволяет утверждать, что при x_0' , достаточно близком к x^* , последовательность (1.3.3) будет сходиться к x^* и для погрешности ε_n' будет выполняться равенство $\varepsilon'_{n+1} = \left(1 - \frac{1}{m}\right) \varepsilon_n' + o(\varepsilon_n')$, показывающее, что в рассматриваемом случае погрешность ε_n' будет изменяться приблизительно по закону геометрической прогрессии со знаменателем $1 - \frac{1}{m}$.

§ 1.4. УЛУЧШЕНИЕ ИТЕРАЦИОННОГО ПРОЦЕССА ПРИ ПОМОЩИ ПРЕОБРАЗОВАНИЯ ЗАДАННОГО УРАВНЕНИЯ

Напомним, что если заданное уравнение $x = \Phi(x)$ таково, что при $x = x^*$ производные от $\Phi(x)$ до порядка $m-1$ равны нулю: $\Phi^{(j)}(x^*) = 0$ ($j = 1, 2, \dots, m-1$) и $\Phi^{(m)}(x^*) \neq 0$, то погрешности приближений $\varepsilon_n = x^* - x_n$ для простого одношагового процесса итерации $x_{n+1} = \Phi(x_n)$ изменяются по следующему закону:

$$\varepsilon_{n+1} = (-1)^{m-1} \frac{1}{m!} \Phi^{(m)}(x^*) \varepsilon_n^m + o(\varepsilon_n^m)$$

и поэтому можно ожидать весьма быстрой сходимости $x_n \rightarrow x^*$, если только x_0 взято достаточно близко к x^* . Если $\Phi(x)$ этим свойством не обладает и, например, $\Phi'(x^*) \neq 0$, то мы, во-первых, можем наверное гарантировать сходимость не всегда, а лишь при условии $|\Phi'(x^*)| < 1$ и, во-вторых, в этом случае закон изменения погрешности ε_n будет

$$\varepsilon_{n+1} = \Phi'(x^*) \varepsilon_n + o(\varepsilon_n),$$

и сходимость будет не столь быстрой и даже может быть медленной, если $|\Phi'(x^*)|$ имеет значение, близкое к 1.

В связи с этим ставится следующая задача: заменить заданное уравнение $x = \Phi(x)$ другим уравнением $x = \Phi(x)$, удовлетворяющим требованиям:

1) уравнение $x = \Phi(x)$ имеет те же решения, что и $x = \Phi(x)$ (или то же решение x^* , когда речь идет о нахождении одного определенного решения, а не всех решений заданного уравнения);

2) для каждого решения x^* уравнения (или для одного определенного решения) должны выполняться условия $\Phi^{(j)}(x^*) = 0$ ($j = 1, 2, \dots, m-1$).

Число m может задаваться наперед, либо оставаться произвольным, удовлетворяющим условию $m > 1$. Примером такого преобразования может служить правило Стеффенсена, когда от уравнения $x = \Phi(x)$ переходят к уравнению (1.3.4). Мы отнесли это правило к предыдущему параграфу, так как идеи, лежащие в основании правила, легче выясняются в задаче преобразования итерационной последовательности, нежели в задаче преобразования уравнений.

Для дальнейшего нам удобнее перейти от уравнения $x = \Phi(x)$ к уравнению $f(x) = \Phi(x) - x = 0$ и пользоваться при преобразованиях функцией $f(x)$. Такой переход связан с тем, что при $x = x^*$ функция $f(x)$ обращается в нуль, а при x , близком к решению x^* , $f(x)$ будет малой величиной. Это дает возможность при нахождении $\Phi(x)$ пользоваться удобным аппаратом степенных разложений. Переход от функции $f(x)$ к $\Phi(x)$ можно рассматривать как некоторый оператор $\Phi = A(f)$, для которого множество допущенных нами к исследованию функций f будет областью определения A и множество функций $\Phi(x)$ — областью значений оператора $A(f)$. Условия перехода от заданного уравнения к новому, указанные выше, налагают слабые ограничения на оператор A и оставляют большой произвол в его выборе, т. е. в выборе соответствия $f \rightarrow \Phi$. Известно несколько таких преобразований. Мы остановимся на двух из них.

1. Будем искать Φ в форме сложной функции

$$\Phi(x) = F[x, f(x)]. \quad (1.4.1)$$

Здесь оператор преобразования $\Phi = A(f)$ определяется выбором функции $F(x, y)$ двух аргументов x и y . Избранная форма представления Φ сужает класс допустимых операторов A , но, как мы увидим ниже, оставляет достаточно большой произвол, чтобы можно было удовлетворить указанным выше двум требованиям, которым мы намерены подчинить $\Phi(x)$. В последующих рассуждениях мы будем предполагать функцию f имеющую непрерывные производные до порядка m , при этом, ради простоты, будем считать первую производную $f'(x)$ отличной от нуля в окрестности решения x^* уравнения $F(x, y)$ предположим имеющей непрерывные частные производные $F_{x^p y^q}$ ($p, q = 0, 1, \dots, m$).

Каждой функции $F(x, y)$, ввиду равенства (1.4.1), отвечает оператор, переводящий m -кратно непрерывно дифференцируемую функцию f в m -кратно непрерывно дифференцируемую функцию $\Phi(x)$. Выясним теперь условия, которым нужно подчинить выбор F , чтобы Φ , отвечающая заданной нам произвольной, но фиксированной функции f , обладала нужными свойствами. В вычислениях удобно воспользоваться степенным разложением F . Применим формулу Тейлора и разложим $F(x, y)$ по аргументу y :

$$\Phi(x) = F(x, f(x)) = a_0(x) + a_1(x)f(x) + a_2(x)f^2(x) + \dots + a_{m-1}(x)f^{m-1}(x) + f^m(x)R_m(x, f(x)), \quad (1.4.2)$$

$$a_p(x) = \frac{1}{p!} F_{yp}(x, 0), \quad R_m(x, f(x)) = \frac{1}{(m-1)!} \int_0^1 F_{ym}(x, f(x)t) (1-t)^{m-1} dt.$$

Когда x есть решение уравнения $f(x)=0$, правая часть равенства приведет к свободному члену $a_0(x)$ и, если мы хотим, чтобы такое значение x было одновременно решением уравнения $\Phi(x)=x$, мы должны потребовать, чтобы $a_0(x)=x$.

Теперь потребуем, чтобы производные $\Phi^{(p)}(x)$ ($p=1, 2, \dots, m-1$), когда x есть решение уравнения $f(x)=0$, обращались в нуль. Заметим попутно, что производные до порядка $m-1$ от последнего члена правой части $f^m R_m$ равны нулю при $f=0$.

$$\left. \begin{aligned} &\Phi'(x) |_{f=0} = 1 + a_1(x)f'(x) = 0, \\ &\Phi''(x) |_{f=0} = 2a_1'(x)f'(x) + a_1(x)f''(x) + 2a_2(x)f'^2(x) = 0, \\ &\Phi'''(x) |_{f=0} = a_1(x)f'''(x) + 3a_1'(x)f''(x) + 3a_1''(x)f'(x) + \\ &\quad + 6a_2'(x)f'^2(x) + 6a_2(x)f'(x)f''(x) + 6a_3(x)f'^3(x) = 0, \\ &\dots\dots\dots \\ &\Phi^{(m-1)}(x) |_{f=0} = (m-1)a_1^{(m-2)}(x)f'(x) + \dots + a_1(x)f^{(m-1)}(x) + \dots + \\ &\quad + (m-1)!a_{m-1}(x)[f'(x)]^{m-1} = 0. \end{aligned} \right\} \quad (1.4.3)$$

Из полученной системы последовательно могут быть найдены $a_1(x), a_2(x), \dots, a_{m-1}(x)$. После этого остается произвольным остаточный член $f^m R_m$, или, что равносильно, $F_{y^m}(x, y)$. Наиболее просто — отбросить $f^m R_m$, т. е. положить $F_{y^m} \equiv 0$. Тогда для $\Phi(x)$ получим следующее выражение:

$$\Phi(x) = x + a_1(x)f(x) + a_2(x)f^2(x) + \dots + a_{m-1}(x)f^{m-1}(x), \quad (1.4.4)$$

где $a_p(x)$ ($p=1, 2, \dots, m-1$) определяется из системы (1.4.3).

При $m=2$

$$a_1(x) = -\frac{1}{f'(x)}, \quad \Phi(x) = x - \frac{f(x)}{f'(x)}$$

и соответствующий итерационный процесс будет

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}.$$

Он совпадает, как мы увидим ниже, с процессом Ньютона, рассмотренным в § 1.7. При $m=3$ из системы (1.4.3) найдем

$$a_1 = -\frac{1}{f'(x)}, \quad a_2 = -\frac{f''(x)}{2f'^3(x)}, \quad \Phi(x) = x - \frac{f(x)}{f'(x)} - \frac{f''(x)f^2(x)}{2f'^3(x)}.$$

Итерационный процесс имеет вид

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} - \frac{f''(x_n)f^2(x_n)}{2f'^3(x_n)}. \quad (1.4.4a)$$

Для $n=4$ получим

$$\Phi(x) = x - \frac{f(x)}{f'(x)} - \frac{f''(x)f^2(x)}{2f'^3(x)} - \frac{f^3(x)}{12} - \frac{3f''^2(x) - f'''(x)f'(x)}{f'^5(x)}$$

и итерационный процесс

$$\begin{aligned} x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} - \frac{f^2(x_n)}{2} \cdot \frac{f''(x_n)}{[f'(x_n)]^3} - \frac{f^3(x_n)}{12} \times \\ \times \frac{3[f''(x_n)]^2 - f'''(x_n)f'(x_n)}{[f'(x_n)]^5}. \end{aligned} \quad (1.4.4b)$$

В общем случае, когда m имеет любое значение, итерационные приближения вычисляются по правилу

$$x_{n+1} = \Phi(x_n),$$

где функция Φ определена равенством (1.4.4).

Остановимся на характеристике закона изменения погрешности $\varepsilon_n = x^* - x_n$. Для этой цели можно было бы воспользоваться либо не вполне определенным равенством (1.2.4), содержащим недостаточно точно известную величину $o(\varepsilon^m)$, либо более определенным равенством (1.2.7), позволяющим точнее оценить величину ε_n . Оба эти

равенства, примененные к нашей задаче, содержали бы производную $\Phi^{(m)}$. Для произвольного значения m функция $\Phi(x)$ сложно выражается через f и вычисление производной от нее порядка m , если не воспользоваться упрощающими соображениями, является затруднительным. Мы укажем сейчас сравнительно простое выражение для $\Phi^{(m)}$ через функцию f и ей обратную.

Пусть производная $f'(x)$ отлична от нуля на некотором отрезке $[a, b]$, содержащем внутри себя решение x^* . Так как $f'(x)$ сохраняет знак на $[a, b]$, $f(x)$ будет там монотонной функцией и будет иметь обратную функцию $x=g(y)$. Последняя будет определена на некотором отрезке $[c, d]$, являющемся областью значений $f(x)$ для $x \in [a, b]$. $g(y)$ будет иметь m непрерывных производных, как и f .

Очевидно, $x^*=g(0)$. Выберем произвольную точку $y \in [c, d]$ и построим разложение x^* по формуле Тейлора, как это сделано ниже:

$$x^*=g(0)=g(y-y)=g(y)+\sum_{i=1}^{m-1}(-1)^i \frac{g^{(i)}(y)}{i!} y^i+(-1)^m \frac{g^{(m)}(\eta)}{m!} y^m,$$

где η лежит между 0 и y , или, если подставить сюда вместо y значение $f(x)$ и заметить, что $g(y)=x$,

$$\begin{aligned} x^*&=x+\sum_{i=1}^{m-1}(-1)^i \frac{F^{(i)}[f(x)]}{i!} f^i(x)+(-1)^m \frac{F^{(m)}[f(\xi)]}{m!} f^m(x)= \\ &=x+\sum_{i=1}^{m-1} b_i(x) f^i(x)+(-1)^m \frac{F^{(m)}[f(\xi)]}{m!} f^m(x)=\Psi(x)+R(x). \end{aligned} \quad (1.4.5)$$

Теперь легко проверить, что $\Psi(x)$ совпадает с $\Phi(x)$. В самом деле, $\Psi(x)$ является многочленом степени m от $f(x)$, подобно $\Phi(x)$, но, может быть, лишь с другими коэффициентами.

При $f=0$ $\Psi(x)=x$. Далее, вычислим производную порядка k ($k=1, 2, \dots, m-1$) от обеих частей последнего равенства при $f=0$. Левая часть равенства x^* не зависит от x и $[x^*]^{(k)}=0$. Так как остаток $R(x)$ содержит множителем f^m , производная от него при $f=0$ обратится в нуль. Поэтому

$$\Psi^{(k)}(x)|_{f=0}=0 \quad (k=1, 2, \dots, m-1)$$

и сумма $\Psi(x)$ удовлетворяет тем же условиям (1) и (2), каким была подчинена сумма $\Phi(x)$. Но эти условия, как мы видели, определяют единственным образом коэффициенты $a_i(x)$ суммы $\Phi(x)$. Поэтому $a_k(x)=b_k(x)$ ($k=1, 2, \dots, m-1$) и, следовательно, $\Phi(x)=\Psi(x)$.

Если положить в равенстве (1.4.5) $x=x_n$ и вычесть почленно из $x_{n+1}=\Phi(x_n)$, получим

$$x^*-x_{n+1}=R(x_n)=(-1)^m \frac{F^{(m)}[f(\xi)]}{m!} f^m(x_n).$$

Но

$$f(x_n)=-[f(x^*)-f(x_n)]=-(x^*-x_n)f'(\xi), \quad \xi \in [x_n, x^*],$$

поэтому

$$\varepsilon_{n+1}=\frac{F^{(m)}[f(\xi)]}{m!} [f'(\xi)]^m \varepsilon_n^m. \quad (1.4.6)$$

Это равенство и было нашей целью. Заметим попутно, что его можно было бы записать в форме, не содержащей неизвестных величин ξ и ζ , если воспользоваться интегральной формой остатка $R(x)$ и теоремой Лагранжа.

Обозначим

$$\max_{[a, b]} \left| \frac{F^{(m)}[f(\xi)]}{m!} [f'(\zeta)]^m \right| = q.$$

Из (1.4.6) получится рекурсионная оценка для погрешности ε_n : $|\varepsilon_{n+1}| \leq q |\varepsilon_n|^m$. Повторное применение ее дает возможность найти приводимую ниже оценку погрешности

$$|\varepsilon_n| \leq q^{\frac{1}{m-1}} |q^{\frac{1}{m-1}} \varepsilon_0|^m,$$

показывающую, что при условии $|q^{\frac{1}{m-1}} \varepsilon_0| < 1$ ε_n будет очень быстро стремиться к нулю.

2. Возвратимся к уравнению $x = \Phi(x)$. Второй способ преобразования относится к уравнениям более частного вида, когда $\Phi(x)$ есть некоторый многочлен от x степени $k > 1$. Как и выше, нам удобнее перейти к уравнению $f(x) = 0$, где $f(x) = \Phi(x) - x$. Мы будем предполагать, что корни $f(x)$ все являются простыми. Такое предположение не ограничивает задачу, так как понизить кратности корней до единицы мы можем, например, при помощи алгоритма Евклида.

Возьмем некоторый многочлен $P(x)$, выбором которого займемся позже, и составим уравнение

$$x = \Phi(x), \quad \Phi(x) = x - P(x)f(x). \quad (1.4.7)$$

Если мы хотим сделать простой одношаговый итерационный процесс для этого уравнения сходящимся более быстро, чем геометрическая прогрессия, мы должны $P(x)$ избрать так, чтобы на каждом решении $x = x^*$ уравнения $f(x) = 0$ было $\Phi'(x^*) = 0$:

$$1 - P'(x^*)f(x^*) - P(x^*)f'(x^*) = 1 - P(x^*)f'(x^*) = 0.$$

Для нахождения $P(x)$ могут быть указаны алгоритмы, требующие выполнения только арифметических операций.*)

Ввиду простоты корней $f(x)$, многочлены $f(x)$ и $f'(x)$ являются взаимно простыми. Применим к ним алгоритм Евклида. В его записи под $r_s(x)$ ($s=0, 1, \dots$) подразумеваются многочлены, старшие коэффициенты которых приведены к 1. Мы положим также $f'(x) = k_0 r_0(x)$.

$$\begin{aligned} f(x) &= r_0(x)q_0(x) + k_1 r_1(x), \\ r_0(x) &= r_1(x)q_1(x) + k_2 r_2(x), \\ &\dots \\ r_{i-3}(x) &= r_{i-2}(x)q_{i-2}(x) + k_{i-1} r_{i-1}(x), \\ r_{i-2}(x) &= r_{i-1}(x)q_{i-1}(x) + k_i \cdot 1. \end{aligned}$$

*) Нужный нам результат является частным случаем хорошо известной в алгебре многочленов теоремы: если $P(x)$ и $Q(x)$ — взаимно простые многочлены степеней p и q соответственно, то существуют многочлены $M(x)$ и $N(x)$ степеней не выше $q-1$ и $p-1$ такие, что верно равенство $M(x)P(x) + N(x)Q(x) = 1$.

Если из этих равенств последовательно исключить $r_1(x)$, $r_2(x)$, ..., $r_{i-1}(x)$, получится соотношение вида

$$P(x)f'(x) + L(x)f(x) = 1.$$

Простой подсчет степеней множителей $P(x)$ и $L(x)$ покажет, что P и L есть многочлены, степени которых не больше $n-1$ и $n-2$ соответственно.

Положив в последнем равенстве $x=x^*$ и учитывая, что $f(x^*)=0$, получим результат $P(x^*)f'(x^*)=1$, убеждающий нас в том, что многочлен $P(x)$ является искомым.

Нахождение $P(x)$ может быть сведено к решению системы n линейных уравнений. Будем для упрощения записи считать, что коэффициент при высшей степени в $f(x)$ приведен к 1.

Станем делить последовательно произведения $x^{i-1}f'(x)$ ($i=1, 2, \dots, n$) на $f(x)$:

$$\left. \begin{aligned} 1 \cdot f'(x) &= 0 \cdot f(x) + p_1(x), \\ x f'(x) &= q_1(x)f(x) + p_2(x), \\ &\dots \dots \dots \\ x^{n-1}f'(x) &= q_{n-1}(x)f(x) + p_n(x). \end{aligned} \right\} \quad (1.4.8)$$

Выберем числа C_0, C_1, \dots, C_{n-1} так, чтобы было верно тождество

$$C_{n-1}p_1(x) + C_{n-2}p_2(x) + \dots + C_0p_n(x) = 1. \quad (1.4.9)$$

Умножая после этого равенства (1.4.8) последовательно на $C_{n-1}, C_{n-2}, \dots, C_0$ и складывая, получим

$$P(x)f'(x) = M(x)f(x) + 1, \quad (1.4.10)$$

где

$$P(x) = C_0x^{n-1} + C_1x^{n-2} + \dots + C_{n-1}.$$

Положив здесь $x=x^*$, мы получим $P(x^*)f'(x^*)=1$ и вновь убедимся в том, что многочлен $P(x)$ является искомым.

Условие (1.4.9), которому подчинен выбор C_0, C_1, \dots, C_{n-1} , даст для нахождения их систему n линейных уравнений. Если многочлен $p_i(x)$ записать в форме $p_i(x) = p_{i0}x^{n-1} + p_{i1}x^{n-2} + \dots + p_{in}$, система будет следующей:

$$\left. \begin{aligned} p_{10}C_{n-1} + p_{20}C_{n-2} + \dots + p_{n0}C_0 &= 0, \\ p_{11}C_{n-1} + p_{21}C_{n-2} + \dots + p_{n1}C_0 &= 0, \\ \dots \dots \dots \\ p_{1n-2}C_{n-1} + p_{2n-2}C_{n-2} + \dots + p_{nn-2}C_0 &= 0, \\ p_{1n-1}C_{n-1} + p_{2n-1}C_{n-2} + \dots + p_{nn-1}C_0 &= 1. \end{aligned} \right\} \quad (1.4.11)$$

Отметим, наконец, что для нахождения многочленов $p_i(x)$ ($i=1, 2, \dots, n$) не обязательно составлять равенства (1.4.8), так как многочлены $p_i(x)$ могут быть найдены последовательно из рекурсионного соотношения

$$xp_i(x) = p_{i0}f(x) + p_{i+1}(x), \quad p_1(x) = f'(x), \quad (1.4.12)$$

которое сразу же получится, если равенство $x^{i-1}f'(x) = m_{i-1}(x)f(x) + r_i(x)$ умножить на x и из произведения $xr_i(x)$ выделить часть $p_{i0}f(x)$, делящуюся нацело на $f(x)$. Остаток после выделения $xr_i(x) - p_{i0}f(x)$ будет многочленом степени меньшей n и должен совпадать с $p_{i+1}(x)$.

§ 1.5. ПОНЯТИЕ ОБ ОБЩЕЙ ТЕОРИИ МЕТОДА ИТЕРАЦИИ. ТЕОРЕМА О СЖАТЫХ ОТОБРАЖЕНИЯХ

Пусть X есть произвольное множество элементов x . Допустим, что на X определен оператор $y = \varphi(x)$, значения которого принадлежат тому же множеству: $\varphi(x) \in X$.

Рассмотрим уравнение

$$x = \varphi(x). \quad (1.5.1)$$

Ему можно придать наглядное содержание. Оператор $y = \varphi(x)$ каждому элементу $x \in X$ ставит в соответствие некоторый элемент y из X . Часто говорят, что оператор φ отображает множество X в себя.

Уравнение (1.5.1) означает, что нужно найти элементы x множества X , которые при отображении φ не изменяются. Такие элементы называются неподвижными. К разысканию их применим следующий итерационный алгоритм: элемент x_0 считается заданным и последующие приближения определяются правилом

$$x_{n+1} = \varphi(x_n).$$

Для него могут быть поставлены вопросы об осуществимости построения и сходимости последовательности x_n , подобно тому, как эти вопросы ставились для аналогичного правила (1.2.2) в случае одного численного уравнения.

Метод итерации применим к весьма широкому классу уравнений, и в истории математики известны многочисленные случаи полезных его приложений не только к теории уравнений, но и в вычислительной практике. Теория метода к настоящему времени доведена до большой общности. Особенно большие успехи в этом направлении были достигнуты за последние два десятилетия, когда был применен в исследованиях абстрактный аппарат функционального анализа.*) Мы не будем излагать теорию метода во всей общности и ограничимся тем, что докажем лишь одну из простых теорем о сходимости итерационной последовательности, вполне сходную с теоремой 1 § 1.2.

Множество X будем считать полным метрическим пространством с метрикой $\rho(x, y)$.

Теорема 1. Пусть выполняются условия:

1) оператор $\varphi(x)$ определен в замкнутой шаровой окрестности S начального элемента x_0 :

$$\rho(x, x_0) \leq \delta; \quad (1.5.2)$$

2) для любых двух элементов x и y из S выполняется неравенство

*) Необходимые для чтения настоящего параграфа сведения из функционального анализа можно найти в добавлении I.

$$\rho[\varphi(x), \varphi(y)] \leq q\rho(x, y) \quad (0 \leq q < 1), \quad (1.5.3)$$

где q не зависит от x и y ;

3) для начального элемента верно неравенство

$$\rho[\varphi(x_0), x_0] \leq m;$$

4) числа δ, q, m подчинены условию

$$\frac{m}{1-q} \leq \delta. \quad (1.5.4)$$

Тогда:

1) приближение x_n , вычисляемое по правилу $x_{n+1} = \varphi(x_n)$, может быть построено для любого n и x_n принадлежит области S ;

2) последовательность x_n сходится к некоторому элементу из S :

$$\lim x_n = x^* \quad (x^* \in S);$$

3) предельный элемент x^* есть решение заданного уравнения:

$$x^* = \varphi(x^*);$$

4) для x_n верно неравенство

$$\rho(x_n, x^*) \leq \frac{m}{1-q} q^n.$$

Доказательство. Покажем сначала, что приближения x_n ($n=1, 2, \dots$) могут быть построены, принадлежат области S и для приближений смежных номеров верно неравенство

$$\rho(x_n, x_{n+1}) \leq m q^n. \quad (1.5.5)$$

Проверим это для $n=0$. Так как x_0 принадлежит S , $x_1 = \varphi(x_0)$ может быть построено. Далее, по предположению (3),

$$\rho(x_0, x_1) = \rho[x_0, \varphi(x_0)] \leq m$$

и неравенство (1.5.5) для $n=0$ выполняется.

Допустим теперь, что x_0, x_1, \dots, x_n построены, принадлежат S и для них выполняются неравенства $\rho(x_k, x_{k+1}) \leq m q^k$ ($k=0, 1, \dots, n-1$). Так как $x_n \in S$ и оператор φ на элементе x_n определен, $x_{n+1} = \varphi(x_n)$ может быть построено. Далее, ввиду предположения (2),

$$\rho(x_n, x_{n+1}) \leq q\rho(x_{n-1}, x_n).$$

По индуктивному допущению

$$\rho(x_{n-1}, x_n) \leq mq^{n-1}$$

и поэтому будет

$$\rho(x_n, x_{n+1}) \leq mq^n.$$

Наконец, если применить к расстоянию $\rho(x_0, x_{n+1})$ несколько раз аксиому треугольника и неравенство $\rho(x_k, x_{k+1}) \leq mq^k$ ($k=0, 1, \dots, n$), получим

$$\begin{aligned} \rho(x_0, x_{n+1}) &\leq \rho(x_0, x_1) + \rho(x_1, x_2) + \dots + \\ &+ \rho(x_n, x_{n+1}) \leq m + mq + \dots + mq^n \leq \frac{m}{1-q} \leq \delta. \end{aligned}$$

Элемент x_{n+1} принадлежит, следовательно, области S .

Покажем теперь, что для последовательности x_n выполняется признак Больцано — Коши.

$$\begin{aligned} \rho(x_n, x_{n+p}) &\leq \rho(x_n, x_{n+1}) + \rho(x_{n+1}, x_{n+2}) + \dots + \rho(x_{n+p-1}, x_{n+p}) \leq \\ &\leq mq^n + mq^{n+1} + \dots + mq^{n+p-1} = \frac{m}{1-q} (q^n - q^{n+p}) \leq \frac{m}{1-q} q^n. \end{aligned}$$

Так как $0 \leq q < 1$, величина $\frac{m}{1-q} q^n$ при больших n будет меньше заданного $\varepsilon > 0$ и признак Больцано — Коши действительно выполняется.

Ввиду полноты пространства X существует элемент x^* , к которому сходится x_n :

$$x_n \rightarrow x^* \quad (n \rightarrow \infty).$$

Легко проверить, что x^* принадлежит S . Действительно, если в неравенстве $\rho(x_0, x_n) \leq \delta$ неограниченно увеличивать n , то, ввиду непрерывной зависимости расстояния $\rho(x_0, x_n)$ от x_n и $x_n \rightarrow x^*$, в пределе получится $\rho(x_0, x^*) \leq \delta$ и, следовательно, x^* принадлежит S .

Далее, $\rho[\varphi(x_n), \varphi(x^*)] \leq q\rho(x_n, x^*)$, и так как $\rho(x_n, x^*) \rightarrow 0$, то $\varphi(x_n) \rightarrow \varphi(x^*)$. Если же заметить, что $\varphi(x_n) = x_{n+1} \rightarrow x^*$, в пределе получим $x^* = \varphi(x^*)$ и элемент x^* есть решение заданного уравнения.

Осталось проверить еще утверждение (4) теоремы о скорости сходимости.

Несколькими строками выше было получено неравенство $\rho(x_n, x_{n+p}) \leq \frac{m}{1-q} q^n$. Если здесь неограниченно увеличивать p и принять во внимание, что при этом $x_{n+p} \rightarrow x^*$ и $\rho(x_n, x_{n+p}) \rightarrow \rho(x_n, x^*)$, в пределе получится $\rho(x_n, x^*) \leq \frac{m}{1-q} q^n$. Этим заканчивается доказательство теоремы 1.

Дополним теорему 1 еще утверждением о единственности решения. Заметим предварительно, что в условии (2) теоремы мы считали коэффициент сжатия $\frac{\rho[\varphi(x), \varphi(y)]}{\rho(x, y)}$ ограниченным числом q , не зависящим от элементов x и y . Для доказательства единственности решения такое предположение о равномерной ограниченности коэффициента сжатия не нужно и достаточно считать $\rho[\varphi(x), \varphi(y)] < \rho(x, y)$.

Теорема 2. На всяком множестве элементов, где для любых двух элементов x, y выполняется неравенство $\rho[\varphi(x), \varphi(y)] < \rho(x, y)$, уравнение $x = \varphi(x)$ может иметь не больше одного решения.

Доказательство. Пусть уравнение имеет два различных решения: $x=\varphi(x)$ и $y=\varphi(y)$ и $\rho(x, y) > 0$.

Ho

$$\rho(x, y) = \rho[\varphi(x), \varphi(y)] < \rho(x, y),$$

что невозможно, и допущение $\rho(x, y) > 0$ является неверным.

§ 1.6. МЕТОД ИТЕРАЦИИ ДЛЯ СИСТЕМ УРАВНЕНИЙ

Пусть дана система n численных уравнений с n неизвестными x_1, x_2, \dots, x_n . Применение метода итерации к ней требует, как и в случае одного уравнения, приведения системы к каноническому виду. Мы будем предполагать, что система дана в нужной форме

$$\left. \begin{aligned} x_1 &= \varphi_1(x_1, x_2, \dots, x_n), \\ x_2 &= \varphi_2(x_1, x_2, \dots, x_n), \\ . &. \\ x_n &= \varphi_n(x_1, x_2, \dots, x_n). \end{aligned} \right\} \quad (1.6.1)$$

С целью сделать более компактной запись, рассмотрим n -мерное числовое пространство R_n , элементами которого являются упорядоченные совокупности n произвольных чисел. Для истолкования в R_n системы (1.6.1) мы должны взять два элемента, первый из которых будет служить для изображения совокупности (x_1, x_2, \dots, x_n) значений аргументов — его мы обозначим x — и второй, который мы обозначим φ , будет служить для изображения соответствующих значений функций $(\varphi_1, \varphi_2, \dots, \varphi_n)$. Элемент φ будет некоторой функцией элемента x . Зависимость $y = \varphi(x)$ можно рассматривать как отображение R_n или части R_n в R_n . Заданная система запишется в короткой форме

$$x = \varphi(x), \quad (1.6.2)$$

и решение ее равносильно нахождению таких элементов x в R_n , которые при отображении $y = \varphi(x)$ переходят в себя.

Допустим, что нами избран исходный элемент для итерации $x^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})$. Все следующие приближения к точному решению $x^* = (x_1^*, x_2^*, \dots, x_n^*)$ находятся по правилу

$$x^{(k+1)} = \varphi(x^{(k)}). \quad (1.6.3)$$

Прежде чем рассматривать теоремы о сходимости $x^{(k)}$ к x^* , мы остановимся на выяснении при помощи нестрогих рассуждений наглядной картины поведения $x^{(n)}$ вблизи точного решения x^* .

Рассмотрим погрешность $x^* - x^{(k)} = \varepsilon^{(k)}$. Соотношение между $\varepsilon^{(k+1)}$ и $\varepsilon^{(k)}$ получится, если в (1.6.3) вместо $x^{(k)}$ и $x^{(k+1)}$ подставить их значения $x^{(k)} = x^* - \varepsilon^{(k)}$, $x^{(k+1)} = x^* - \varepsilon^{(k+1)}$:

$$x^* - \varepsilon^{(k+1)} = \varphi(x^* - \varepsilon^{(k)}).$$

Это равенство элементов R_n равносильно n численным равенствам

$$x_i^* - \varepsilon_i^{(k+1)} = \varphi_i(x_1^* - \varepsilon_1^{(k)}, \dots, x_n^* - \varepsilon_n^{(k)}) \quad (i=1, 2, \dots, n).$$

Предположим теперь, что $x^{(k)}$ лежит близко к x^* и погрешности $\varepsilon_i^{(k)}$ являются малыми величинами. Разложим правую часть по степеням $\varepsilon_i^{(k)}$ и выделим из разложения линейную часть. Если принять во внимание, что $x_i^* = \varphi_i(x_1^*, x_2^*, \dots, x_n^*)$, то предыдущее равенство приведет к следующему соотношению между $\varepsilon^{(n)}$ и $\varepsilon^{(n+1)}$:

$$\varepsilon_i^{(k+1)} = \sum_{j=1}^n \frac{\partial}{\partial x_j} \varphi_i(x_1^*, \dots, x_n^*) \varepsilon_j^{(k)} + o(\max_j |\varepsilon_j^{(k)}|).$$

Отсюда видно, что при выполнении одной итерации погрешность $\varepsilon^{(k)} = (\varepsilon_1^{(k)}, \dots, \varepsilon_n^{(k)})$ претерпевает приблизительно линейное преобразование

$$\varepsilon^{(k+1)} \approx A \varepsilon^{(k)}. \quad (1.6.4)$$

В этой записи под $\varepsilon^{(k)}$ подразумевается n -мерный вектор с координатами $(\varepsilon_1^{(k)}, \dots, \varepsilon_n^{(k)})$; A есть значение на точном решении x^* матрицы Якоби системы функций φ_i :

$$A = \begin{bmatrix} \frac{\partial}{\partial x_1} \varphi_1(x_1^*, \dots, x_n^*) & \dots & \frac{\partial}{\partial x_n} \varphi_1(x_1^*, \dots, x_n^*) \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1} \varphi_n(x_1^*, \dots, x_n^*) & \dots & \frac{\partial}{\partial x_n} \varphi_n(x_1^*, \dots, x_n^*) \end{bmatrix}.$$

Чтобы сделать наглядным изменение, которое претерпевают погрешности $\varepsilon_i^{(n)}$ при преобразовании, мы выполним линейную замену переменных. Для простоты будем считать, что все элементарные делители матрицы A являются простыми. Существует такая неособенная матрица S , что A представима в виде

$$A = S^{-1}[\lambda_1, \lambda_2, \dots, \lambda_n] S.$$

Здесь $\lambda_1, \dots, \lambda_n$ — собственные значения матрицы A и $[\lambda_1, \dots, \lambda_n]$ — диагональная матрица, элементы которой, расположенные на диагонали, указаны в скобках.

$$\varepsilon^{(n+1)} \approx S^{-1}[\lambda_1, \lambda_2, \dots, \lambda_n] S \varepsilon^{(n)}. \quad (1.6.5)$$

Примем произведение $S \varepsilon^{(n)}$ за новый вектор, подлежащий исследованию, и положим $\eta^{(n)} = S \varepsilon^{(n)}$. Так как такое преобразование — неособенное, стремление $\varepsilon^{(n)}$ к нулю при $n \rightarrow \infty$ равносильно сходимости к нулю $\eta^{(n)}$.

Умножив обе части равенства (1.6.5) слева на S , получим равенство, дающее закон изменения $\eta^{(n)}$ при одном шаге итерационного процесса:

$$\eta^{(n+1)} \approx [\lambda_1, \lambda_2, \dots, \lambda_n] \eta^{(n)},$$

что равносильно n численным равенствам

$$\eta_i^{(k+1)} \approx \lambda_i \eta_i^{(k)} \quad (i=1, 2, \dots, n). \quad (1.6.6)$$

Каждая из величин $\eta_i^{(k)}$ ($i=1, \dots, n$) с изменением k будет изменяться приблизительно как геометрическая прогрессия со своим знаменателем λ_i . Если все λ_i имеют модули, меньшие единицы:

$$|\lambda_i| < 1 \quad (i=1, 2, \dots, n),$$

то весьма вероятно, что $\eta_i^{(k)} \rightarrow 0$ при $k \rightarrow \infty$. Вместе с $\eta_i^{(k)}$ будут стремиться к нулю все погрешности $\varepsilon_i^{(k)}$ и итерационная последовательность $x^{(k)}$ будет сходиться к решению x^* .

Если же среди λ_i будут большие или равные единице по модулю, то нельзя гарантировать стремление к нулю всех $\eta_i^{(k)}$ и, следовательно, погрешностей $\varepsilon_i^{(k)}$ при $k \rightarrow \infty$. Более того, если среди λ_i будут числа, строго большие единицы по модулю, то при малых $\eta_i^{(k)}$, когда равенства (1.6.6) будут иметь малые относительные погрешности, некоторые $\eta_i^{(k+1)}$, по видимому, будут иметь модули большие, нежели $\eta_i^{(k)}$, и поэтому решение x^* будет «элементом отталкивания» для итерационных приближений $x^{(k)}$.

Изложенное выше описание поведения погрешностей $\epsilon_i^{(k)}$ является очень наглядным, но его еще недостаточно для составления удовлетворительной картины поведения приближений $x^{(k)}$ вблизи x^* , так как описание имеет преимущественно качественный характер и лишено количественных оценок. Главным недостатком описания является следующий факт: оно основано на приближенном равенстве (1.6.4), которое получено из предшествующего точного равенства путем отбрасывания «малого слагаемого» $o(\max_j |\epsilon_j^{(k)}|)$. Мы не оценивали это слагаемое и не оценивали

окрестность x^* , в которой оно не исказит существенно картины поведения $x^{(n)}$, полученной нами на основании рассмотрения приближенного равенства (1.6.4). В изложении мы ограничились лишь указанной нами качественной картиной и не занимались количественной стороной вопроса.

Перейдем теперь к теоремам о сходимости $x^{(n)}$ к x^* . Нужные нам результаты мы получим как следствие из теоремы 1 § 1.5 о сжатых отображениях. В теореме имеются в виду метрические пространства, и, чтобы воспользоваться ей в нашей задаче, мы должны ввести в n -мерное числовое пространство метрику.

В одномерном числовом пространстве, геометрически изображаемом точками числовой оси, естественной метрикой является абсолютное значение разности между числами: $\rho(x, y) = |x - y|$. Оно равно длине отрезка между точками x и y на оси. В многомерном пространстве нет единственной естественной метрики, и в разных задачах бывает целесообразно пользоваться различными определениями метрики в зависимости от условий задачи и целей, которые мы преследуем.

Мы будем пользоваться тремя определениями расстояния, которые наиболее часто употребляются в практике вычислений.

1. Кубическая метрика. Ее мы обозначим $\rho_m(x, y)$ и определим равенством

$$\rho_m(x, y) = \max_i |x_i - y_i|.$$

2. Октаэдрическая метрика. Обозначается $\rho_s(x, y)$ и определяется так:

$$\rho_s(x, y) = \sum_{i=1}^n |x_i - y_i|.$$

3. Шаровая метрика. Обозначим ее $\rho_l(x, y)$ и определим равенством

$$\rho_l^2(x, y) = \sum_{i=1}^n (x_i - y_i)^2.$$

Нам нужно будет для формулировки теорем подсчитать оценку коэффициента растяжения при преобразовании $y = \varphi(x)$ в каждой из этих

метрик. Функции φ_i будем предполагать непрерывно дифференцируемыми.

1. Случай кубической метрики. В области $\rho_m(x, x^{(0)}) \leq \delta$ возьмем два произвольных элемента $x = (x_1, x_2, \dots, x_n)$ и $y = (y_1, y_2, \dots, y_n)$ и рассмотрим координаты $\varphi_i(x)$ и $\varphi_i(y)$ их изображений.

$$\begin{aligned} |\varphi_i(x) - \varphi_i(y)| &= |\varphi_i(x_1, \dots, x_n) - \varphi_i(y_1, \dots, y_n)| = \\ &= \left| \sum_{j=1}^n \left(\frac{\partial \varphi_i}{\partial x_j} \right)^* (x_j - y_j) \right| \leq \sum_{j=1}^n \left| \left(\frac{\partial \varphi_i}{\partial x_j} \right)^* \right| \max_j |x_j - y_j| = \\ &= \left(\sum_{j=1}^n \left| \frac{\partial \varphi_i}{\partial x_j} \right| \right)^* \rho_m(x, y). \end{aligned}$$

Под символом $(\)^*$ здесь понимается значение функции, стоящей в скобках, в некоторой точке прямолинейного отрезка, соединяющего x и y . Положение этой точки зависит от x, y, i . Чтобы получить оценку, не зависящую от x, y и индекса i , заменим

$$\left(\sum_{j=1}^n \left| \frac{\partial \varphi_i}{\partial x_j} \right| \right)^* \quad \text{на} \quad \max_i \max_x \sum_{j=1}^n \left| \frac{\partial \varphi_i}{\partial x_j} \right|,$$

где \max_x означает наибольшее значение в области $\rho_m(x, x^{(0)}) \leq \delta$. После этого мы получим

$$\rho_m[\varphi(x), \varphi(y)] \leq \max_i \max_x \sum_{j=1}^n \left| \frac{\partial \varphi_i}{\partial x_j} \right| \rho_m(x, y).$$

Отсюда видно, что в качестве оценки для коэффициента растяжения $\frac{\rho_m[\varphi(x), \varphi(y)]}{\rho_m(x, y)}$ может быть взята величина

$$q = \max_i \max_x \sum_{j=1}^n \left| \frac{\partial \varphi_i}{\partial x_j} \right|. \quad (1.6.7)$$

2. Случай октаэдрической метрики. Рассматривается область $\rho_s(x, x^{(0)}) \leq \delta$. Вычисления, сходные с теми, которые приведены в предыдущем случае, покажут, что верно неравенство

$$\rho_s[\varphi(x), \varphi(y)] = \sum_{i=1}^n |\varphi_i(x) - \varphi_i(y)| \leq q \rho_s(x, y), \quad (1.6.8)$$

$$q = \sum_{i=1}^n \max_j \max_x \left| \frac{\partial \varphi_i}{\partial x_j} \right|.$$

Под $\max F$ подразумевается наибольшее значение функции F в области $\rho_s(x, x^{(0)}) \leq \delta$.

3. Случай шаровой метрики. Рассматриваемая область есть

$$\rho_l(x, y) = \left\{ \sum_{i=1}^n (x_i - y_i)^2 \right\}^{\frac{1}{2}} \leq \delta.$$

Верны следующие оценки:

$$\begin{aligned} |\varphi_i(x) - \varphi_i(y)|^2 &= \left| \sum_{j=1}^n \left(\frac{\partial \varphi_i}{\partial x_j} \right)^* (x_j - y_j) \right|^2 \leq \\ &\leq \sum_{j=1}^n \left[\left(\frac{\partial \varphi_i}{\partial x_j} \right)^* \right]^2 \sum_{j=1}^n (x_j - y_j)^2 \leq \max_x \sum_{j=1}^n \left(\frac{\partial \varphi_i}{\partial x_j} \right)^2 \rho_l^2(x, y). \\ \rho_l^2[\varphi(x), \varphi(y)] &= \sum_{i=1}^n [\varphi_i(x) - \varphi_i(y)]^2 \leq q^2 \rho_l^2(x, y), \quad (1.6.9) \end{aligned}$$

$$q^2 = \sum_{i=1}^n \max_x \sum_{j=1}^n \left(\frac{\partial \varphi_i}{\partial x_j} \right)^2.$$

Укажем сейчас теоремы, дающие в трех метриках условия, достаточные для сходимости итеративной последовательности. Ввиду того, что теоремы читаются во всех трех случаях аналогично, мы приведем полную формулировку лишь для кубической метрики, для двух же других укажем только изменения, которые должны быть внесены в теорему.

Теорема 1. Пусть выполняются условия:

1) функции $\varphi_i(x_1, \dots, x_n)$ ($i=1, 2, \dots, n$) определены и непрерывно дифференцируемы в области

$$|x_i - x_i^{(0)}| \leq \delta \quad (i=1, 2, \dots, n); \quad (1.6.10_1)$$

2) удовлетворяют там неравенствам

$$\max_x \sum_{j=1}^n \left| \frac{\partial \varphi_i}{\partial x_j} \right| \leq q < 1 \quad (i=1, 2, \dots, n); \quad (1.6.11_1)$$

3) для начальных приближений $x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)}$ выполняются условия

$$|x_i^{(0)} - \varphi_i(x_1^{(0)}, \dots, x_n^{(0)})| \leq m \quad (i=1, 2, \dots, n); \quad (1.6.12_1)$$

4) для чисел δ, q и m соблюдается неравенство

$$\frac{m}{1-q} \leq \delta.$$

Тогда:

1) система (1.6.1) в области (1.6.10₁) имеет решение $x^* = (x_1^*, \dots, x_n^*)$, к которому сходится итерационная последовательность приближений $x^{(h)} = (x_1^{(h)}, \dots, x_n^{(h)})$, вычисляемая по правилу (1.6.3);

2) скорость сходимости может быть охарактеризована неравенством

$$|x_i^* - x_i^{(h)}| \leq \frac{m}{1-q} q^h \quad (i=1, 2, \dots, n). \quad (1.6.13_1)$$

Заметим также, что из теоремы 2 предыдущего параграфа вытекает, что система (1.6.1) может иметь в области (1.6.10₁) не более одного решения, если там выполняется условие

$$\max_x \sum_{j=1}^n \left| \frac{\partial \varphi_i}{\partial x_j} \right| < 1 \quad (i=1, \dots, n). \quad (1.6.14_1)$$

Для октаэдрической метрики теорема и замечание к ней читаются так же, если условия (1.6.10₁), (1.6.11₁), (1.6.12₁) и неравенства (1.6.13₁), (1.6.14₁) заменить соответственно на следующие:

$$\sum_{i=1}^n |x_i - x_i^{(0)}| \leq \delta, \quad (1.6.10_2)$$

$$\sum_{i=1}^n \max_j \max_x \left| \frac{\partial \varphi_i}{\partial x_j} \right| \leq q < 1, \quad (1.6.11_2)$$

$$\sum_{i=1}^n |x_i^{(0)} - \varphi_i(x_1^{(0)}, \dots, x_n^{(0)})| \leq m, \quad (1.6.12_2)$$

$$\sum_{i=1}^n |x_i^* - x_i^{(k)}| \leq \frac{m}{1-q} q^k, \quad (1.6.13_2)$$

$$\sum_{i=1}^n \max_j \max_x \left| \frac{\partial \varphi_i}{\partial x_j} \right| < 1. \quad (1.6.14_2)$$

Наконец, в случае шаровой метрики неравенства (1.6.10₁)—(1.6.14₁) должны быть заменены на указанные ниже:

$$\sum_{i=1}^n (x_i - x_i^{(0)})^2 \leq \delta^2, \quad (1.6.10_3)$$

$$\sum_{i=1}^n \max_x \sum_{j=1}^n \left(\frac{\partial \varphi_i}{\partial x_j} \right)^2 \leq q < 1, \quad (1.6.11_3)$$

$$\sum_{i=1}^n [x_i^{(0)} - \varphi_i(x_1^{(0)}, \dots, x_n^{(0)})]^2 \leq m^2, \quad (1.6.12_3)$$

$$\left\{ \sum_{i=1}^n (x_i^* - x_i^{(k)})^2 \right\}^{1/2} \leq \frac{m}{1-q} q^k, \quad (1.6.13_3)$$

$$\sum_{i=1}^n \max_x \sum_{j=1}^n \left(\frac{\partial \varphi_i}{\partial x_j} \right)^2 < 1. \quad (1.6.14_3)$$

§ 1.7. МЕТОД НЬЮТОНА. СЛУЧАЙ ОДНОГО ЧИСЛЕННОГО УРАВНЕНИЯ

Метод Ньютона является весьма общим и применим к решению широкого класса нелинейных операторных уравнений, в частности нелинейных численных уравнений. Его значение заключается в том, что он позволяет решение нелинейного уравнения свести к решению последовательности линейных задач. Делается это при помощи выделения из заданного нелинейного уравнения главной линейной части. С формальной точки зрения метод Ньютона можно рассматривать как частный случай метода итерации, но, как мы увидим ниже, он основан на идее, совершенно отличной от идеи, лежащей в основе метода итерации.

Мы ознакомимся с идеей метода Ньютона на примере его применения к решению одного уравнения с одной численной неизвестной величиной.

Пусть дано нелинейное уравнение $f(x) = 0$, где x есть численная переменная и f — достаточно гладкая функция. Обозначим x^* точное решение уравнения. Предположим, что каким-либо путем нами указано для x^* исходное приближение x_0 , и поставим перед собой задачу построить алго-

ритм для его уточнения при помощи построения линейного уравнения, приближенно заменяющего заданное и являющегося его главной частью. Нам удобнее от точного решения x^* перейти к другой величине, которую в условиях проблемы можно считать малой. Рассмотрим погрешность исходного приближения $\varepsilon = x^* - x_0$. Нахождение x^* или погрешности ε — равносильные задачи, но для наших целей удобнее пользоваться ε , так как, как правило, x_0 бывает близким к x^* и ε — малой величиной.

Для составления уравнения, из которого может быть найдена погрешность, достаточно в равенство $f(x^*) = 0$ вместо x^* подставить его значение $x_0 + \varepsilon$:

$$f(x_0 + \varepsilon) = 0. \quad (1.7.1)$$

С целью выделения отсюда главной линейной части можно разложить $f(x_0 + \varepsilon)$ по формуле Тейлора, ограничиваясь в ней лишь членами, линейными относительно ε , и относя все члены, содержащие ε в степенях выше первой, в остаток:

$$f(x_0) + \varepsilon f'(x_0) + o(\varepsilon) = 0.$$

Считая ε малой величиной, отбросим остаток разложения $o(\varepsilon)$. Мы получим после этого линейное уравнение для погрешности ε , близкое к (1.7.1) и отличающееся от него только на малую величину $o(\varepsilon)$ выше первого порядка:

$$f(x_0) + \varepsilon f'(x_0) \approx 0. \quad (1.7.2)$$

Решая его, мы найдем не точное значение погрешности ε , а лишь приближенное ее значение, которое обозначим ε_0 . Численная величина

$\varepsilon_0 = -\frac{f(x_0)}{f'(x_0)}$, как можно было ожидать, будет главной частью ε . Добавляя ее к x_0 , получим исправленное приближение $x_1 = x_0 + \varepsilon_0$, и можно надеяться, что оно будет ближе к x^* , нежели x_0 :

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Его можно в свою очередь улучшить, пользуясь теми же соображениями, и т. д. В результате получим последовательность приближений x_n ($n=0, 1, \dots$), в которой каждое следующее значение x_{n+1} находится по правилу

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad (n=0, 1, 2, \dots). \quad (1.7.3)$$

Условием возможности процесса приближений (1.7.3) является выполнение двух требований:

- 1) все x_n принадлежат области задания $f(x)$;
- 2) $f'(x_n) \neq 0$.

Процесс Ньютона совпадает с простым одношаговым итерационным процессом для уравнения $x = \varphi(x)$ при $\varphi(x) = x - \frac{f(x)}{f'(x)}$; с которым мы встречались в § 1.4, где занимались вопросом уточнения метода итерации.

Правило (1.7.3) имеет простой геометрический смысл. В плоскости с системой координат xOy построим график l функции $y = f(x)$ (рис. 1.7.1). Абсцисса точки пересечения l с осью Ox является решением x^* уравнения $f(x) = 0$. Отметим на оси Ox точку x_n и рассмотрим на графике l соответствующую точку $M_n(x_n, f(x_n))$. В этой точке проведем касательную T_n к линии l и найдем точку пересечения T_n с осью Ox . Уравнение касательной есть $y - f(x_n) = f'(x_n)(x - x_n)$. Положив здесь $y = 0$, мы получим уравнение для нахождения абсциссы пересечения T_n и Ox . Обозначим эту абсциссу x_{n+1} . Для ее нахождения получится уравнение $-f(x_n) = f'(x_n)(x_{n+1} - x_n)$ и $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$, что совпадает

с x_{n+1} в правиле Ньютона (1.7.3). Поэтому правило Ньютона геометрически означает, что решение x^* уравнения находится приближенно при помощи замены линии l касательной прямой T_n .

Теперь постараемся путем нестрогих рассуждений выяснить наглядную картину поведения приближений x_n вблизи точного решения x^* . Нам удобнее перейти к погрешности $\varepsilon_n = x^* - x_n$. Если из равенства $x^* = x^*$ почленно вычесть (1.7.3), получим следующее соотношение между ε_n и ε_{n+1} :

$$\varepsilon_{n+1} = -\frac{f(x_n) + \varepsilon_n f'(x_n)}{f'(x_n)}.$$

Заметив, что

$$0 = f(x^*) = f(x_n + \varepsilon_n) = f(x_n) + \varepsilon_n f'(x_n) + \frac{1}{2} \varepsilon_n^2 f''(x_n + \Theta \varepsilon_n),$$

мы легко получим выражение для ε_{n+1} , дающее описание закона изменения порядка погрешности на одном шаге приближений:

$$\varepsilon_{n+1} = -\frac{1}{2} \cdot \frac{f''(x_n + \Theta \varepsilon_n)}{f'(x_n)} \varepsilon_n^2.$$

Отсюда можно без труда получить одну из простейших теорем, в которой указываются достаточные условия сходимости к x^* ньютоновой последовательности x_n . Мы применим это равенство к более частному

вопросу. При малых ε_n множитель у ε_n^2 будет мало отличаться от $-\frac{1}{2} \cdot \frac{f''(x^*)}{f'(x^*)}$ и

$$\varepsilon_{n+1} = -\frac{1}{2} \cdot \frac{f''(x^*)}{f'(x^*)} \varepsilon_n^2 + o(\varepsilon_n^2).$$

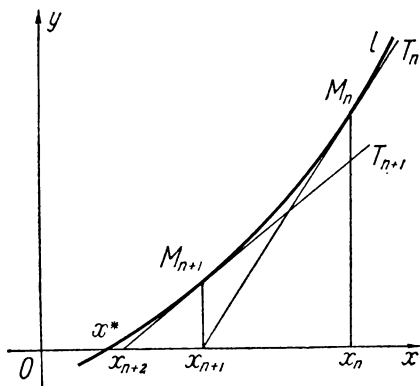


Рис. 1.7.1

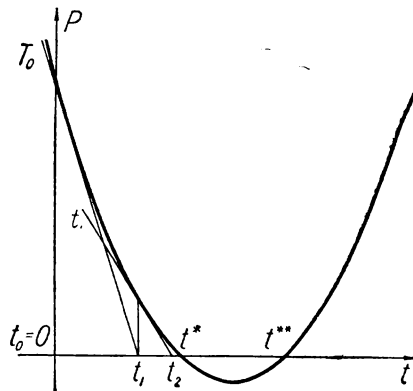


Рис. 1.7.2

Пренебрегая здесь малой высшего порядка малости, найдем нужный нам приближенный закон изменения погрешности ε_n :

$$\varepsilon_{n+1} \approx -\frac{1}{2} \cdot \frac{f''(x^*)}{f'(x^*)} \varepsilon_n^2 = -\alpha \varepsilon_n^2. \quad (1.7.4)$$

Из него следует, что $\frac{\varepsilon_{n+1}}{\varepsilon_n} \approx \left(\frac{\varepsilon_n}{\varepsilon_{n-1}} \right)^2$ и отношение погрешностей при одном шаге будет возводиться в квадрат.

Предположим, что равенства вида (1.7.4)

$$\varepsilon_i \approx -\alpha \varepsilon_{i-1}^2 \quad (i=1, 2, \dots, n)$$

выполняются с достаточной точностью. С их помощью нетрудно выразить ε_n через ε_0 . Получим

$$\varepsilon_n \approx -(\alpha \varepsilon_0)^{2^{n-1}} \varepsilon_0 = -\left[\frac{1}{2} \cdot \frac{f''(x^*)}{f'(x^*)} \varepsilon_0 \right]^{2^{n-1}} \varepsilon_0, \quad (1.7.5)$$

При $|\alpha\epsilon_0| < 1$ погрешность ϵ_n будет быстро стремиться к нулю и процесс Ньютона (1.7.3) будет весьма быстро сходиться к точному решению.

Перейдем теперь к теореме о разрешимости уравнения и сходимости к решению ньютоновой последовательности приближений. Для теоремы, которую мы докажем, особую роль играет вопрос о сходимости метода Ньютона для уравнения второй степени

$$P(t) = at^2 + bt + c = 0,$$

где a, b, c — вещественные числа и $b^2 - 4ac \geq 0$. Корни его вещественные. Меньший из них обозначим t^* и больший — t^{**} . За исходное значение t_0

примем любое значение $t_0 \neq -\frac{b}{2a}$. Геометрически ясно, что если t_0 лежит вне отрезка $[t^*, t^{**}]$, то ньютоновская последовательность t_n будет монотонно сходиться к ближайшему корню многочлена. Если же t_0 лежит внутри $[t^*, t^{**}]$, то уже первый шаг вычислений выведет из этого отрезка, t_1 окажется лежащим вне $[t^*, t^{**}]$ и после этого получится указанная выше монотонная последовательность t_n . Чтобы убедиться в правильности сказанного, достаточно рассмотреть рис. 1.7.2.

Теорема 1. Пусть выполняются условия:

1) функция $f(x)$ определена и дважды непрерывно дифференцируема на отрезке

$$|x - x_0| \leq \delta, \quad (\alpha)$$

при этом $|f''(x)| \leq K$ при всяких x на этом отрезке;

$$2) f'(x_0) \neq 0 \text{ и } \frac{1}{|f'(x_0)|} \leq B;$$

3) выполняется неравенство

$$\left| \frac{f(x_0)}{f'(x_0)} \right| \leq \eta;$$

4) для B, K, η соблюдено условие

$$h = BK\eta \leq \frac{1}{2};$$

5) верно неравенство

$$\frac{1 - \sqrt{1 - 2h}}{h} \eta \leq \delta.$$

Тогда:

- 1) последовательность $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$ ($n=0, 1, 2, \dots$) может быть построена и является сходящейся: $x_n \rightarrow x^*$;
- 2) предельное значение x^* есть решение уравнения $f(x) = 0$;
- 3) верна оценка скорости сходимости

$$|x^* - x_n| \leq t^* - t_n, \quad (1.7.6)$$

где t_n — ньютонова последовательность приближений $t_{n+1} = t_n - \frac{P(t_n)}{P'(t_n)}$ к меньшему корню t^* уравнения $P(t) = \frac{1}{2}Kt^2 - \frac{1}{B}t + \frac{\eta}{B} = 0$, построенная при $t_0 = 0$.

Доказательство. Имеем: $t^* = \frac{1 - \sqrt{1-2h}}{h}\eta$. Так как $0 < h \leq \frac{1}{2}$, корень t^* есть действительное и, очевидно, неотрицательное число. Последовательность t_n будет монотонно возрастающей и сходящейся к t^* .

Покажем при помощи индукции, что все x_n ($n=0, 1, \dots$) могут быть найдены, лежат внутри области (α) и для них имеет место оценка

$$|x_{n+1} - x_n| \leq t_{n+1} - t_n. \quad (1.7.7)$$

Сначала проверим это для $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$. Так как x_0 лежит внутри отрезка (α) и $f'(x_0) \neq 0$, приближение x_1 может быть найдено. Заметим, что дробь

$$\frac{1 - \sqrt{1-2h}}{h} = \frac{2}{1 + \sqrt{1-2h}}$$

при $0 < h \leq \frac{1}{2}$ будет изменяться в границах $(1, 2]$. Поэтому

$$|x_1 - x_0| = \left| \frac{f(x_0)}{f'(x_0)} \right| \leq \eta < \frac{1 - \sqrt{1-2h}}{h} \eta \leq \delta$$

и x_1 лежит внутри (α) . По условию $t_0 = 0$. Так как

$$t_1 = t_0 - \frac{P(t_0)}{P'(t_0)} = \frac{\frac{\eta}{B}}{\frac{1}{B}} = \eta, \quad |t_1 - t_0| = \eta \quad \text{и} \quad |x_1 - x_0| \leq \eta,$$

то неравенство (1.7.7) для x_1 и x_0 выполнено.

Допустим теперь, что x_0, x_1, \dots, x_n могут быть построены, лежат внутри (α) и для них выполняются неравенства

$$|x_{k+1} - x_k| \leq t_{k+1} - t_k \quad (k=0, 1, \dots, n-1).$$

По предположению x_n лежит внутри (α) и $f(x_n), f'(x_n)$ имеют смысл. Для установления возможности построения x_{n+1} достаточно проверить, что $f'(x_n) \neq 0$. Но

$$\begin{aligned} |f'(x_n)| &= \left| f'(x_0) + \int_{x_0}^{x_n} f''(t) dt \right| \geq \frac{1}{B} - K |x_n - x_0| \geq \\ &\geq \frac{1}{B} - K |(x_n - x_{n-1}) + (x_{n-1} - x_{n-2}) + \dots + (x_1 - x_0)| \geq \\ &\geq \frac{1}{B} - K [(t_n - t_{n-1}) + (t_{n-1} - t_{n-2}) + \dots + (t_1 - t_0)] = \\ &= \frac{1}{B} - K(t_n - t_0) = \frac{1}{B} - Kt_n = -P'(t_n). \end{aligned}$$

Рис. 1.7.2 показывает, что $-P'(t_n) > 0$ и, стало быть, $|f'(x_n)| > 0$. Оценим теперь $f(x_n)$. Воспользуемся теїлоровыми разложениями около точки x_{n-1}

$$f(x_n) = f(x_{n-1}) + (x_n - x_{n-1})f'(x_{n-1}) + \int_{x_{n-1}}^{x_n} f''(t)(x_n - t) dt.$$

Ввиду $x_n = x_{n-1} + \frac{f(x_{n-1})}{f'(x_{n-1})}$ сумма неинтегральных членов в правой части равна нулю и

$$f(x_n) = \int_{x_{n-1}}^{x_n} f''(t)(x_n - t) dt.$$

Отсюда получается оценка

$$|f(x_n)| \leq K \frac{(x_n - x_{n-1})^2}{2}$$

или, так как по индуктивному допущению $|x_n - x_{n-1}| < t_n - t_{n-1}$,

$$|f(x_n)| \leq \frac{1}{2} K(t_n - t_{n-1})^2.$$

Аналогичные вычисления для t_n и $P(t_n)$ дадут

$$P(t_n) = \int_{t_{n-1}}^{t_n} P''(t) (t_n - t) dt = \frac{1}{2} K(t_n - t_{n-1})^2.$$

Поэтому $|f(x_n)| \leq P(t_n)$. Значит,

$$|x_{n+1} - x_n| = \left| \frac{f(x_n)}{f'(x_n)} \right| \leq -\frac{P(t_n)}{P'(t_n)} = t_{n+1} - t_n$$

и неравенство (1.7.7) для x_n и x_{n+1} выполняется. Нам осталось проверить, что x_{n+1} лежит внутри отрезка (α) . Имеем

$$\begin{aligned} |x_{n+1} - x_0| &= |(x_{n+1} - x_n) + (x_n - x_{n-1}) + \dots + (x_1 - x_0)| \leq (t_{n+1} - t_n) + \\ &+ (t_n - t_{n-1}) + \dots + (t_1 - t_0) = t_{n+1} - t_0 = t_{n+1} < t^* = \frac{1 - \sqrt{1 - 2h}}{h} \eta \leq \delta \end{aligned}$$

и x_{n+1} , следовательно, лежит внутри (α) .

Докажем сходимость x_n . Для этого достаточно проверить выполнение признака Больцано — Коши.

$$\begin{aligned} |x_{n+p} - x_n| &= |(x_{n+p} - x_{n+p-1}) + (x_{n+p-1} - x_{n+p-2}) + \dots + (x_{n+1} - x_n)| \leq \\ &\leq (t_{n+p} - t_{n+p-1}) + (t_{n+p-1} - t_{n+p-2}) + \dots + (t_{n+1} - t_n) = t_{n+p} - t_n. \end{aligned}$$

Но последовательность t_n сходится и для нее признак Больцано — Коши выполняется. Полученное же неравенство доказывает, что он будет выполняться и для x_n .

Обозначим $\lim_{n \rightarrow \infty} x_n = x^*$. Утверждение о скорости сходимости (1.7.6)

получится сразу же, если перейти к пределу при $p \rightarrow \infty$ в неравенстве $|x_{n+p} - x_n| \leq t_{n+p} - t_n$, так как $x_{n+p} \rightarrow x^*$ и $t_{n+p} \rightarrow t^*$.

Наконец, если в правиле Ньютона $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$ считать, что $n \rightarrow \infty$, мы ввиду $x_n \rightarrow x^*$ и $x_{n+1} \rightarrow x^*$ получим $\frac{f(x_n)}{f'(x_n)} \rightarrow 0$ и, так как $f'(x_n)$ является величиной, ограниченной, отсюда вытекает $f(x_n) \rightarrow 0$, что, по непрерывности f в точке x^* , влечет за собой $f(x^*) = 0$. Предел x^* действительно есть решение заданного уравнения.

Оценка скорости сходимости (1.7.6) в условиях теоремы 1 является неулучшаемой, так как она достигается для квадратного уравнения

$$\frac{1}{2} K x^2 - \frac{1}{B} x + \frac{\eta}{B} = 0 \text{ при } x_0 = 0. \text{ Но она недостаточно наглядна ввиду}$$

того, что требует знания разности $t^* - t_n$. Изменение ее при росте n может быть охарактеризовано численной таблицей. В квадратном уравнении содержится два существенных параметра. С целью освободиться от одного из них введем новую переменную τ , положив $t = \eta\tau$. Подстановка в $P(t)$ даст

$$P(\eta\tau) = \frac{\eta}{B} \left(\frac{1}{2} h\tau^2 - \tau + 1 \right).$$

Рассмотрим теперь квадратное уравнение

$$\varphi(\tau) = \frac{1}{2} h\tau^2 - \tau + 1 = 0 \quad (1.7.8)$$

и построим для него ньютоновскую последовательность приближений $\tau_0 = 0$, $\tau_{n+1} = \tau_n - \frac{\varphi(\tau_n)}{\varphi'(\tau_n)}$ ($n = 0, 1, \dots$). Меньший корень квадратного уравнения есть $\tau^* = \frac{1 - \sqrt{1 - 2h}}{h}$, и к нему будет сходиться последовательность τ_n . Легко можно убедиться в том, что $t_n = \eta\tau_n$.

Уравнение (1.7.8) содержит только один параметр h , и это упрощает задачу табулирования. Для разности $\tau^* - \tau_n$ может быть составлена таблица значений в зависимости от h и от n . Краткая таблица такого рода [6, стр. 58], позволяющая составить представление о скорости убывания $\tau^* - \tau_n$ с ростом n , приведена ниже.

Таблица значений $\tau^* - \tau_n$

$h \backslash n$	0	1	2	3	4	5
0,05	1,026	$0,263 \cdot 10^{-1}$	$0,183 \cdot 10^{-4}$	$0,877 \cdot 10^{-11}$	$0,203 \cdot 10^{-23}$	
0,10	1,056	$0,557 \cdot 10^{-1}$	$0,173 \cdot 10^{-3}$	$0,166 \cdot 10^{-8}$	$0,155 \cdot 10^{-18}$	
0,15	1,089	$0,889 \cdot 10^{-1}$	$0,698 \cdot 10^{-3}$	$0,436 \cdot 10^{-7}$	$0,171 \cdot 10^{-15}$	
0,20	1,127	0,127	$0,202 \cdot 10^{-2}$	$0,525 \cdot 10^{-6}$	$0,356 \cdot 10^{-13}$	
0,25	1,172	0,172	$0,491 \cdot 10^{-2}$	$0,425 \cdot 10^{-5}$	$0,319 \cdot 10^{-11}$	$0,180 \cdot 10^{-23}$
0,30	1,225	0,225	$0,109 \cdot 10^{-1}$	$0,278 \cdot 10^{-4}$	$0,184 \cdot 10^{-9}$	$0,802 \cdot 10^{-20}$
0,35	1,292	0,292	$0,230 \cdot 10^{-1}$	$0,166 \cdot 10^{-3}$	$0,885 \cdot 10^{-8}$	$0,250 \cdot 10^{-16}$
0,40	1,382	0,382	$0,486 \cdot 10^{-1}$	$0,101 \cdot 10^{-2}$	$0,459 \cdot 10^{-6}$	$0,942 \cdot 10^{-13}$
0,45	1,519	0,519	0,110	$0,749 \cdot 10^{-2}$	$0,395 \cdot 10^{-4}$	$0,111 \cdot 10^{-8}$
0,50	2	1	0,5	0,25	0,125	$0,625 \cdot 10^{-1}$

Можно для $x^* - x_n$ получить более грубую, но более наглядную, чем (1.7.6) оценку.

Теорема 2. При соблюдении условий теоремы 1, для разности $x^* - x_n$ верна оценка

$$|x^* - x_n| \leq \frac{1}{2^{n-1}} (2h)^{2^{n-1}} \eta. \quad (1.7.9)$$

Доказательство. Согласно формуле Тейлора, если учесть, что по определению τ_n

$$\varphi(\tau_{n-1}) + (\tau_n - \tau_{n-1})\varphi'(\tau_{n-1}) = 0,$$

то

$$\begin{aligned} \varphi(\tau_n) &= \varphi(\tau_{n-1}) + (\tau_n - \tau_{n-1})\varphi'(\tau_{n-1}) + \frac{1}{2} \varphi''(\xi) (\tau_n - \tau_{n-1})^2 = \\ &= \frac{1}{2} h (\tau_n - \tau_{n-1})^2, \\ \varphi'(\tau_n) &= h\tau_n - 1, \quad \tau_{n+1} - \tau_n = -\frac{\varphi(\tau_n)}{\varphi'(\tau_n)} = \frac{1}{2} \frac{h}{1 - h\tau_n} (\tau_n - \tau_{n-1})^2. \end{aligned} \quad (1.7.10)$$

Так как $\tau_0 = 0$, $\tau_1 = 1$ и $h \leq \frac{1}{2}$, при $n=1$ получим

$$\tau_2 - \tau_1 = \frac{1}{2} \frac{h}{1 - h} \leq \frac{1}{2}.$$

Отсюда вытекает оценка для τ_2 :

$$\tau_2 = \tau_1 + (\tau_2 - \tau_1) \leq 1 + \frac{1}{2} = \frac{3}{2}.$$

Для $n=2$ из (1.7.10) следует

$$\begin{aligned} \tau_3 - \tau_2 &= \frac{1}{2} \frac{h}{1 - h\tau_2} (\tau_2 - \tau_1)^2 \leq \frac{1}{2} \frac{h}{1 - \frac{3}{2}h} \frac{1}{2^2} \leq \frac{1}{4}, \\ \tau_3 &= \tau_2 + (\tau_3 - \tau_2) \leq \frac{3}{2} + \frac{1}{4} = \frac{7}{4}. \end{aligned}$$

Продолжая такие оценки, для произвольных значений $n=0, 1, 2, \dots$ найдем

$$\tau_n \leq 2 - 2^{1-n}. \quad (1.7.11)$$

По определению τ_n и ввиду $\varphi(\tau^*) = 0$, имеем

$$\tau^* - \tau_n = \tau^* - \tau_{n-1} + \frac{\varphi(\tau_{n-1})}{\varphi'(\tau_{n-1})} = - \frac{1}{\varphi'(\tau_{n-1})} [\varphi(\tau^*) - \varphi(\tau_{n-1}) - (\tau^* - \tau_{n-1})\varphi'(\tau_{n-1})].$$

По формуле же Тейлора

$$\begin{aligned} \varphi(\tau^*) - \varphi(\tau_{n-1}) - (\tau^* - \tau_{n-1})\varphi'(\tau_{n-1}) &= \frac{1}{2} \varphi''(\xi) (\tau^* - \tau_{n-1})^2 = \\ &= \frac{1}{2} h (\tau^* - \tau_{n-1})^2, \quad \varphi'(\tau_{n-1}) = h\tau_{n-1} - 1 \end{aligned}$$

и, стало быть,

$$\tau^* - \tau_n = \frac{1}{1 - h\tau_{n-1}} \cdot \frac{1}{2} h (\tau^* - \tau_{n-1})^2.$$

На основании (1.7.11)

$$1 - h\tau_{n-1} \geq 1 - \frac{1}{2} (2 - 2^{2-n}) = 2^{1-n},$$

и поэтому

$$\tau^* - \tau_n \leq 2^{n-2} h (\tau^* - \tau_{n-1})^2. \quad (1.7.12)$$

Применим это неравенство последовательно для $n=1, 2, \dots$, приняв во внимание, что $\tau^* = \frac{1 - \sqrt{1-2h}}{h} \leq 2$.

Для $n=1$ находим

$$\tau^* - \tau_1 \leq 2^{-1} h (\tau^* - \tau_0)^2 \leq 2h.$$

Для $n=2$

$$\tau^* - \tau_2 \leq h (\tau^* - \tau_1)^2 \leq h (2h)^2 = \frac{1}{2} (2h)^3.$$

Продолжая такие оценки, получим

$$\tau^* - \tau_n \leq \frac{1}{2^{n-1}} (2h)^{2^n-1}.$$

Этот результат доказывает теорему, так как

$$|x^* - x_n| \leq t^* - t_n = \eta (\tau^* - \tau_n) \leq \frac{1}{2^{n-1}} (2h)^{2^n-1} \eta.$$

Сделаем еще дополнение к теореме 1'. В ней, в частности, было доказано, что при выполнении условий, перечисленных там, уравнение

$f(x)=0$ будет иметь на отрезке $|x-x_0|\leq\delta$ решение. Мы укажем сейчас область, в которой можно гарантировать единственность решения.

Теорема 3. Пусть для функции $f(x)$ выполнены условия теоремы 1 с тем различием, что $\delta\geq\frac{1+\sqrt{1-2h}}{h}$. Если $h<\frac{1}{2}$, то решение x^* уравнения $f(x)=0$ будет единственным в области $|x-x_0|\leq\gamma$, где

$$\gamma<t^{**}=\frac{1+\sqrt{1-2h}}{h}\eta. \quad (1.7.13)$$

Если же $h=\frac{1}{2}$, то решение x^* будет единственным в области $|x-x_0|\leq\gamma$ при

$$\gamma=t^{**}=\frac{1+\sqrt{1-2h}}{h}\eta=2\eta. \quad (1.7.14)$$

Доказательство. Рассмотрим сначала случай $h<\frac{1}{2}$ и допустим, что \bar{x} есть решение рассматриваемого уравнения, принадлежащее области $|x-x_0|\leq\gamma$. Так как для γ выполнено условие (1.7.13), то можно положить

$$|\bar{x}-x_0|\leq\Theta(t^{**}-t_0)=\Theta t^{**}, \quad 0\leq\Theta<1. \quad (1.7.15)$$

Ввиду $f(\bar{x})=0$ верно равенство

$$\bar{x}-x_1=-\frac{1}{f'(x_0)}[f(\bar{x})-f(x_0)-(\bar{x}-x_0)f'(x_0)].$$

По формуле Тейлора

$$f(\bar{x})=f(x_0)+(\bar{x}-x_0)f'(x_0)+\int_{x_0}^{\bar{x}}f''(t)(\bar{x}-t)dt$$

и, следовательно,

$$\bar{x}-x_1=-\frac{1}{f'(x_0)}\int_{x_0}^{\bar{x}}f''(t)(\bar{x}-t)dt.$$

Отсюда

$$|\bar{x}-x_1|\leq\frac{1}{|f'(x_0)|}\cdot\frac{1}{2}K(\bar{x}-x_0)^2.$$

При доказательстве теоремы 1 нами при оценке $f'(x_n)$ было получено неравенство $|f'(x_n)| \geq |P(t_n)|$, где $P(t)$ есть вспомогательный квадратный многочлен. Если, кроме того, воспользоваться неравенством (1.7.15), получим оценку

$$|\bar{x} - x_1| \leq \frac{\Theta^2}{|P'(t_0)|} \cdot \frac{1}{2} K t^{**2}.$$

Наконец, так как

$$\begin{aligned} \frac{1}{P'(t_0)} \cdot \frac{1}{2} K t^{**2} &= \frac{1}{P'(t_0)} [P(t^{**}) - P(t_0) - (t^{**} - t_0) P'(t_0)] = \\ &= -t^{**} + t_0 - \frac{P(t_0)}{P'(t_0)} = t_1 - t^{**}, \end{aligned}$$

найдем нужную нам оценку

$$|\bar{x} - x_1| \leq \Theta^2 (t^{**} - t_1).$$

Сравним ее с (1.7.15). Она получается из (1.7.15) заменой Θ на Θ^2 и x_0, t_0 на x_1, t_1 .

Применив это правило n раз, получим неравенство

$$|\bar{x} - x_n| \leq \Theta^{2^n} (t^{**} - t_n) < \Theta^{2^n} t^{**}. \quad (1.7.16)$$

Так как $\Theta < 1$ и не зависит от n , отсюда следует, что $x_n \rightarrow \bar{x}$, и поэтому $\bar{x} = x^*$.

Когда $h = \frac{1}{2}$, число Θ может равняться единице. Оба корня многочлена $P(t)$ тогда совпадают: $t^{**} = t^*$ и, так как $t_n \rightarrow t^*$ ($n \rightarrow \infty$), из (1.7.16) вытекает, что $|x_n - \bar{x}| \rightarrow 0$. Этим теорема доказана.

§ 1.8. ОБ УТОЧНЕНИЯХ И ИЗМЕНЕНИЯХ МЕТОДА НЬЮТОНА

Метод Ньютона есть метод линеаризации задачи. Он является одним из старейших вычислительных методов решения уравнений. Его история имеет почти трехсотлетнюю давность. Метод очень часто применяется, и поэтому естественно, что было сделано много попыток, преследующих цель изменить его либо в сторону увеличения скорости сходимости, либо — упрощения вычислений.

В большинстве видоизменений нарушалась идея, лежащая в основе метода Ньютона, — идея линейности уравнения, из которого должно быть найдено улучшенное приближение. Очень часто или указывались рекур-

сионные правила для получения следующего приближения, или для их нахождения требовалось решать нелинейное уравнение. Ниже мы дадим очень краткое описание некоторых направлений, в которых можно изменять метод Ньютона. Более полное изложение приведем лишь для небольшого числа видоизменений, где сохранялась линейность уравнения, служащего для получения улучшенного приближения.

Начнем с вопроса увеличения скорости сходимости. Прежде всего напомним, что правило Ньютона $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$ является частным случаем одношагового простого правила итерации $x_{n+1} = \varphi(x_n)$ при $\varphi(x) = x - \frac{f(x)}{f'(x)}$. Поэтому все приемы увеличения скорости сходимости

метода итерации, о которых мы говорили в § 1.4, могут быть перенесены на метод Ньютона для уравнения $f(x) = 0$. К такому виду относятся, например, итерационные правила (1.4.4а) и (1.4.4б), полученные там.

При вычислениях по правилу Ньютона мы составляем таблицы значений трех величин $x_k, f(x_k), f'(x_k)$ ($k=0, 1, \dots$). Пусть таблица доведена до значений $x_n, y_n = f(x_n), y'_n = f'(x_n)$ и вычислению подлежит x_{n+1} . По существу дела мы выполняем линейное интерполирование функции $y = f(x)$ по одному узлу x_n и известным значениям в этой точке функции $y_n = f(x_n)$ и производной от нее $y'_n = f'(x_n)$ и полагаем приближенно $y \approx f(x_n) + f'(x_n)(x - x_n)$. Затем находим то значение x , при котором y обращается в нуль, и принимаем его за x_{n+1} .

Но мы можем надеяться найти x_{n+1} с большей точностью, если выполним интерполирование $y = f(x)$ по значениям f и f' не в одном узле x_n , а в нескольких предшествующих узлах $x_n, x_{n-1}, \dots, x_{n-k}$, или интерполирование обратной функции $x = F(y)$ по значениям F и производной

$\frac{dx}{dy} = F'(y) = \frac{1}{f'(x)}$ в узлах $y_n, y_{n-1}, \dots, y_{n-k}$. В обоих случаях мы за-

меняем на некоторый многочлен, вообще говоря нелинейный, функцию f или ей обратную и отступаем от идеи приближенной линеаризации задачи. Поэтому мы не будем рассматривать указанные методы в параграфах, посвященных методу Ньютона, и более подробное ознакомление с ними отложим до гл. 4.

Можно указать одну задачу в проблеме увеличения точности метода Ньютона, которая решается в рамках идеи линеаризации. До настоящего места мы предполагали, что первая производная $f'(x)$ не обращается в нуль на интервале, где лежат все приближения x_n . В отдельных вопросах мы предполагали, кроме того, $f'(x^*) \neq 0$, т. е. корень x^* считали однократным. Сейчас мы выясним вопрос о скорости сходимости ньютоновой последовательности в случае, когда кратность решения x^* выше первой. Убедимся сначала, что сходимость $x_n \rightarrow x^*$ в этом случае сильно замед-

ляется. Допустим, что решение x^* имеет кратность m и, стало быть, разложение f по формуле Тейлора в окрестности решения x^* имеет форму

$$f(x) = a_m(x-x^*)^m + a_{m+1}(x-x^*)^{m+1} + \dots + a_{m+p}(x-x^*)^{m+p} + R_{m+p}, \quad (1.8.1)$$

$$a_k = \frac{1}{k!} f^{(k)}(x^*) \quad (k=0, 1, \dots).$$

Допустим, что x_n находятся вблизи x^* и погрешности $\varepsilon_n = x^* - x_n$ являются малыми величинами. Правило Ньютона дает следующую связь между погрешностями ε_n и ε_{n+1} :

$$\varepsilon_{n+1} = \varepsilon_n + \frac{f(x^* - \varepsilon_n)}{f'(x^* - \varepsilon_n)}.$$

Сохраняя в разложениях только два главных члена, при помощи (1.8.1) находим:

$$\begin{aligned} f(x^* - \varepsilon_n) &= (-1)^m [a_m \varepsilon_n^m - a_{m+1} \varepsilon_n^{m+1} + \dots], \\ f'(x^* - \varepsilon_n) &= (-1)^{m-1} [m a_m \varepsilon_n^{m-1} - (m+1) a_{m+1} \varepsilon_n^m + \dots], \\ [f'(x^* - \varepsilon_n)]^{-1} &= \frac{(-1)^{m-1}}{m a_m} \varepsilon_n^{m+1} \left[1 + \frac{(m+1) a_{m+1}}{m a_m} \varepsilon_n + \dots \right], \\ \frac{f(x^* - \varepsilon_n)}{f'(x^* - \varepsilon_n)} &= -\frac{\varepsilon_n}{m} \left[1 + \frac{a_{m+1}}{m a_m} \varepsilon_n + \dots \right], \\ \varepsilon_{n+1} &= \left(1 - \frac{1}{m} \right) \varepsilon_n - \frac{a_{m+1}}{m^2 a_m} \varepsilon_n^2 + \dots \end{aligned} \quad (1.8.2)$$

или, если сохранить только один главный член,

$$\varepsilon_{n+1} \approx \left(1 - \frac{1}{m} \right) \varepsilon_n. \quad (1.8.3)$$

Последнее равенство говорит, что ε_n изменяется приблизительно по закону геометрической прогрессии со знаменателем $q = 1 - \frac{1}{m}$, меньшим единицы. Если сравнить этот результат с равенством (1.7.4), дающим закон изменения ε_n в случае простого корня x^* , когда $f'(x^*) \neq 0$, будет видно, что в случае кратного корня сходимость $x_n \rightarrow x^*$ будет намного медленнее.

Можно поставить вопрос о том, как следует изменить вычислительное правило Ньютона, чтобы в случае кратного корня улучшить его сходимость и сделать ее приблизительно такой же быстрой, как и для простого корня. Оказывается, что этого можно достичь весьма просто, оставаясь в границах идеи линейности уравнения для погрешности.

Введем число k , выбор которого совершим позже, и рассмотрим вычислительное правило

$$x_{n+1} = x_n - k \frac{f(x_n)}{f'(x_n)}. \quad (1.8.4)$$

Соответствующий ему закон изменения погрешностей ε_n есть, очевидно,

$$\varepsilon_{n+1} = \left(1 - \frac{k}{m}\right) \varepsilon_n - k \frac{a_{m+1}}{m^2 a_m} \varepsilon_n^2 + \dots$$

Если мы хотим увеличить скорость убывания ε_n при возрастании n , то для этого достаточно положить $k=m$. Тогда главный член в правой части равенства исчезнет и будет верным следующий приближенный закон изменения ε_n :

$$\varepsilon_{n+1} \approx - \frac{a_{m+1}}{m a_m} \varepsilon_n^2 = - \frac{f^{(m+1)}(x^*)}{m(m+1)f^{(m)}(x^*)} \varepsilon_n^2. \quad (1.8.5)$$

Можно сказать поэтому, что если разыскиваемый корень x^* уравнения $f(x)=0$ имеет кратность m , то вычислительное правило

$$x_{n+1} = x_n - m \frac{f(x_n)}{f'(x_n)} \quad (1.8.6)$$

будет иметь приблизительно такой же закон убывания погрешности ε_n для x_n , близких к x^* , как и в основном правиле Ньютона (1.7.3) при $f'(x^*) \neq 0$.

Геометрический смысл правила весьма прост. Будем считать x_n известным и рассмотрим точку $M_n(x_n, f(x_n))$ на графике l функции $f(x)$. Напомним, что в правиле Ньютона (1.7.3) мы проводили в точке M_n касательную к линии l и за x_{n+1} принимали абсциссу точки пересечения касательной с осью Ox . Новое правило (1.8.6) получается из (1.7.3) при помощи замены $f'(x_n)$ на $\frac{1}{m} f'(x_n)$. Последнее означает, что через точку M_n проводится не касательная прямая к l , а вспомогательная прямая, для которой тангенс угла наклона к оси Ox в m раз меньше, чем для касательной прямой. Координата точки пересечения такой прямой с осью Ox и принимается в правиле (1.8.6) за x_{n+1} .

Теперь перейдем к видоизменениям метода Ньютона, имеющим целью упростить вычислительную работу за счет некоторой потери в скорости сходимости последовательности x_n . В методе Ньютона главная доля труда затрачивается на нахождение значений $f(x_n)$ и $f'(x_n)$ и было бы желательно избежать вычисления одной из этих величин вполне или отчасти. В первую очередь естественно отказаться от вычисления значений $f'(x_n)$, так как для суждения о близости x_n к решению знание $f(x_n)$ имеет, вообще говоря, большее значение, чем знание $f'(x_n)$.

1. Начнем с метода секущих. Можно было бы заменить $f'(x_n)$ надлежаще подобранной линейной комбинацией из нескольких значений функции f . Мы рассмотрим наиболее простую из таких замен, когда $f'(x_n)$ приближенно вычисляется по двум последним найденным парам чисел $(x_n, f(x_n))$ и $(x_{n-1}, f(x_{n-1}))$:

$$f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}.$$

Соответствующее правило вычислений следующего приближения будет:

$$x_{n+1} = x_n - \frac{(x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})} = \frac{x_{n-1}f(x_n) - x_nf(x_{n-1})}{f(x_n) - f(x_{n-1})}. \quad (1.8.7)$$

Оно имеет простой геометрический смысл. Возьмем на графике l функции $f(x)$ две точки $M_n[x_n, f(x_n)]$ и $M_{n-1}[x_{n-1}, f(x_{n-1})]$ и проведем через них секущую прямую линию. Уравнение секущей есть

$$\frac{x - x_n}{x_{n-1} - x_n} = \frac{y - f(x_n)}{f(x_{n-1}) - f(x_n)}$$

и правило (1.8.7) означает, что за x_{n+1} принята абсцисса точки пересечения секущей с осью Ox .

По сравнению с основным правилом Ньютона метод секущих имеет особенности, на которые мы хотим обратить внимание. Этот метод является двухшаговым и для начала вычислений требует знания двух исходных приближений x_0 и x_1 к решению x^* .

Условиями возможности осуществления алгоритма (1.8.7) являются:

- 1) принадлежность всех x_n области определения f ;
- 2) выполнение неравенств

$$f(x_n) - f(x_{n-1}) \neq 0 \quad (n = 1, 2, \dots).$$

Остановим свое внимание на случае, когда $f(x_n) - f(x_{n-1}) = 0$. Здесь могут быть две возможности.

1. Пусть $x_{n-1} \neq x_n$. Как видно из равенства

$$x_n = x_{n-1} - \frac{(x_{n-1} - x_{n-2})f(x_{n-1})}{f(x_{n-1}) - f(x_{n-2})},$$

значение $f(x_{n-1}) \neq 0$. Поэтому $f(x_n)$ также отлично от нуля и следующее приближение

$$x_{n+1} = x_n - \frac{(x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})}$$

не может быть построено. Процесс приближений по правилу секущих здесь оборвется на приближении x_n и не приведет к решению.

2. Допустим теперь, что $x_n = x_{n-1}$. Мы считаем при этом, что все приближения $x_0, x_1, \dots, x_{n-1}, x_n$ могут быть построены, x_0, x_1, \dots, x_{n-1} различны между собой и $f(x_{k+1}) - f(x_k) \neq 0$ ($k=0, 1, \dots, n-2$). Из приведенного выше выражения для x_n следует, что $f(x_{n-1}) = 0$ и x_{n-1} есть решение заданного уравнения. В этом случае последовательные приближения осуществимы до значения x_n , при этом два совпадающих последних значения x_{n-1} и x_n являются решениями заданного уравнения.

Изучение правила (1.8.7) начнем с выяснения поведения приближений x_n вблизи решения x^* . Для погрешности $\varepsilon_n = x^* - x_n$ из (1.8.7) получится уравнение

$$\varepsilon_{n+1} = \varepsilon_n + \frac{(\varepsilon_{n-1} - \varepsilon_n)f(x^* - \varepsilon_n)}{f(x^* - \varepsilon_n) - f(x^* - \varepsilon_{n-1})}.$$

Если сюда внести вместо $f(x^* - \varepsilon_n)$ и $f(x^* - \varepsilon_{n-1})$ их разложения по степеням погрешностей

$$f(x^* - \varepsilon_n) = -f'(x^*)\varepsilon_n + \frac{1}{2}f''(x^*)\varepsilon_n^2 + \dots,$$

$$f(x^* - \varepsilon_{n-1}) = -f'(x^*)\varepsilon_{n-1} + \frac{1}{2}f''(x^*)\varepsilon_{n-1}^2 + \dots$$

и выполнить разложение правой части по степеням ε_n и ε_{n-1} , сохраняя лишь члены до второго порядка относительно этих величин, найдем

$$\varepsilon_{n+1} \approx -\frac{f''(x^*)}{2f'(x^*)}\varepsilon_n\varepsilon_{n-1}. \quad (1.8.8)$$

Сравнение полученного равенства с аналогичным равенством (1.7.4), дающим закон изменения погрешности ε_n в основном методе Ньютона, говорит, что погрешности в методе секущих будут изменяться по закону,

близкому к закону изменения ε_n в основном методе, но с несколько меньшей скоростью стремления ε_n к нулю.*)

Приведенная ниже теорема о сходимости последовательности x_n , построенной по методу секущих (1.8.7), высказана в форме, близкой к сходной теореме для основного метода Ньютона. Мы будем предполагать выполненными условия последней теоремы и в соответствии с этим считать известным, что решение x^* существует и является единственным. В теореме для нас полезны в первую очередь сведения о возможности осуществления метода секущих и о сходимости x_n .

Предварительно сделаем два пояснения, полезные для понимания условий теоремы. В методе секущих, как и в основном методе Ньютона, мы должны предполагать, что все последовательные приближения x_n лежат на отрезке, где производная f' не обращается в нуль и f , следовательно, будет монотонно возрастающей или убывающей функцией.

Вычисления начинаются с двух исходных значений, которые были обозначены x_0 и x_1 . Для нахождения x_2 по правилу (1.8.7) безразлично, какое из приближений взять за x_0 и какое за x_1 . Но когда мы перейдем к нахождению x_3 по x_2 и x_1 , становится не безразличным, что принято за x_1 . Рассмотрим значения $f(x_0)$, $f(x_1)$. Они различны по модулю. За x_0 естественно принять то из двух приближений, которому отвечает большее по модулю значение функции и считать $|f(x_0)| > |f(x_1)|$.

Исследование сходимости приближений будет основано на сравнении x_n с приближениями t_n , построенными по методу секущих, к меньшему корню t^* квадратного многочлена

$$P_n(t) = \frac{1}{B} Kt^2 - \frac{1}{B} t + \frac{\eta}{B},$$

известного по теореме 1 § 1.7, когда вычисления начинаются со значений $t_0=0$ и $t_1>0$.

Чтобы сделать возможным сравнение и придать ему простую форму, оказывается, достаточно подчинить x_1 первому неравенству условия (7) теоремы и положить $t_1 = |x_1 - x_0|$.

Кроме того, в доказательстве будет иметь значение приведенное ниже неравенство (1.8.9). Для $n=1$ оно выполняется ввиду выбора t_1 . Для индуктивного доказательства его необходимо, чтобы оно было верным для первого шага вычислений: $|x_2 - x_1| \leq t_2 - t_1$. Ввиду же условий (1) — (6) теоремы, для этого достаточно, как будет выяснено ниже, считать $|f(x_1)| \leq P(t_1)$, что совпадает со вторым неравенством условия (7).

Теорема 1. Пусть для уравнения $f(x)=0$ и исходных приближений x_0 и x_1 выполняются условия:

*) Можно дать более точное сравнение законов убывания ε_n при $n \rightarrow \infty$ для основного метода Ньютона и метода секущих, если решить разностное уравнение (1.8.8). Для наших целей достаточным является качественное заключение, указанное в тексте.

1) $f(x)$ определена и дважды непрерывно дифференцируема на отрезке $|x-x_0| \leq \delta$;

$$2) f'(x_0) \neq 0 \text{ и } \frac{1}{|f'(x_0)|} \leq B;$$

$$3) \left| \frac{f(x_0)}{f'(x_0)} \right| \leq \eta;$$

4) вторая производная $f''(x)$ на отрезке $|x-x_0| \leq \delta$ ограничена по модулю числом K :

$$|f''(x)| \leq K;$$

5) для чисел B, η, K справедливо неравенство

$$h = B\eta K \leq \frac{1}{2};$$

$$6) \delta \geq \frac{1 - \sqrt{1-2h}}{h} \eta;$$

7) для x_1 выполняются неравенства

$$|x_1 - x_0| < \frac{1 - \sqrt{1-2h}}{h} \eta = t^* \text{ и } |f(x_1)| \leq P(|x_1 - x_0|) = P(t_1).$$

Тогда:

1) приближения x_n , определяемые правилом секущих (1.8.7), либо приведут к решению за конечное число шагов, либо могут быть построены при всяком n и образуют сходящуюся последовательность:

$$\lim_{n \rightarrow \infty} x_n = x^*,$$

предельное значение x^* которой есть решение заданного уравнения $f(x) = 0$;

2) скорость сходимости может быть охарактеризована неравенством

$$|x^* - x_n| \leq t^* - t_n, \quad (1.8.9)$$

где t_n — последовательные приближения по методу секущих к меньшему корню t^* квадратного уравнения

$$P(t) = \frac{1}{2} K t^2 - \frac{1}{B} t + \frac{\eta}{B} = 0$$

при $t_0 = 0, t_1 = |x_1 - x_0|$.

Доказательство. Ввиду $0 < h \leq \frac{1}{2}$ корни

$$t^* = \frac{1 - \sqrt{1-2h}}{h} \eta \text{ и } t^{**} = \frac{1 + \sqrt{1-2h}}{h} \eta$$

многочлена $P(t)$ действительны и неотрицательны. Условие

$$t_1 = |x_1 - x_0| < \frac{1 - \sqrt{1-2h}}{h} \eta$$

означает, что t_1 лежит между $t_0=0$ и t^* . Если сделать чертеж, то из него будет видно, что последовательность t_n будет возрастающей и сходящейся к t^* .

Рассмотрим всю последовательность приближений x_n , построенных по правилу секущих (1.8.7). Пусть это будет x_0, x_1, \dots, x_N . Она может быть конечной или бесконечной. В последнем случае мы будем условно считать $N=\infty$. Напомним, что конечной эта последовательность может быть по двум причинам: либо x_N выйдет за границы отрезка $|x-x_0| \leq \delta$, являющегося областью определения f , либо будет $f(x_N) - f(x_{N-1}) = 0$.

Докажем, что для последовательности x_n верны неравенства

$$|x_{n+1} - x_n| \leq t_{n+1} - t_n \quad (n=0, 1, \dots, N-1). \quad (1.8.10)$$

Для этого прибегнем к индукции. При $n=0$ неравенство (1.8.10) выполняется ввиду того, что $t_1 = |x_1 - x_0|$ и $t_0=0$.

Проверим выполнение его для $n=1$. Если воспользоваться формулой Тейлора, получим оценку делителя

$$\begin{aligned} |f(x_1) - f(x_0)| &= |(x_1 - x_0)f'(x_0) + \frac{1}{2}(x_1 - x_0)^2 f''(\xi)| \geq \\ &\geq |x_1 - x_0| \cdot \left[\frac{1}{B} - \frac{1}{2} K |x_1 - x_0| \right] = |x_1 - x_0| \cdot \left[\frac{1}{B} - \frac{1}{2} K t_1 \right]. \end{aligned}$$

Следовательно,

$$|x_2 - x_1| \leq \frac{|f(x_1)|}{\frac{1}{B} - \frac{1}{2} K t_1}.$$

С другой стороны, если проделать такие же вычисления для $t_2 - t_1$, получим ввиду $t_0=0$

$$t_2 - t_1 = - \frac{(t_1 - t_0)P(t_1)}{P(t_1) - P(t_0)} = \frac{P(t_1)}{\frac{1}{B} - \frac{1}{2} K t_1},$$

и так как по условию (7) $|f(x_1)| \leq P(t_1)$, из полученных оценок следует $|x_2 - x_1| \leq t_2 - t_1$.

Допустим теперь, что (1.8.10) выполняется для x_0, x_1, \dots, x_n ($n < N$), и проверим его выполнение для пары x_n, x_{n+1} .

Начнем с оценки $f(x_n)$. По правилу (1.8.7)

$$x_n - x_{n-1} = - \frac{f(x_{n-1})(x_{n-1} - x_{n-2})}{f(x_{n-1}) - f(x_{n-2})}$$

или

$$(x_n - x_{n-2})f(x_{n-1}) - (x_n - x_{n-1})f(x_{n-2}) = 0.$$

Отсюда следует, что

$$L(f, x_n) = \frac{x_n - x_{n-2}}{x_{n-1} - x_{n-2}} f(x_{n-1}) + \frac{x_n - x_{n-1}}{x_{n-2} - x_{n-1}} f(x_{n-2}) = 0.$$

Но $L(f, x_n)$ есть значение в точке x_n линейного многочлена, интерполирующего функцию $f(x)$ по двум ее значениям $f(x_{n-2})$ и $f(x_{n-1})$ в узлах x_{n-2}, x_{n-1} и $f(x_n) - L(f, x_n)$ есть значение остатка интерполирования. По известной теореме об остатке, будем иметь

$$f(x_n) - L(f, x_n) = \frac{(x_n - x_{n-1})(x_n - x_{n-2})}{2} f''(\xi).$$

Здесь ξ есть некоторая средняя точка на отрезке, содержащем x_{n-2}, x_{n-1}, x_n .

На основании формулы Тейлора

$$f(x_n) - f(x_{n-1}) = (x_n - x_{n-1})f'(x_{n-1}) + \frac{1}{2} (x_n - x_{n-1})^2 f''(\eta),$$

где η лежит на отрезке $[x_{n-1}, x_n]$. С помощью последних двух формул получаем

$$x_{n+1} - x_n = - \frac{(x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})} = - \frac{\frac{1}{2} (x_n - x_{n-1})(x_n - x_{n-2})f''(\xi)}{f'(x_{n-1}) + \frac{1}{2} (x_n - x_{n-1})f''(\eta)}. \quad (1.8.11)$$

По индуктивному предположению,

$$|x_n - x_{n-1}| \leq t_n - t_{n-1}$$

и

$$|x_n - x_{n-2}| \leq |x_n - x_{n-1}| + |x_{n-1} - x_{n-2}| \leq (t_n - t_{n-1}) + (t_{n-1} - t_{n-2}) = t_n - t_{n-2}.$$

Числитель в последней части (1.8.11) может быть оценен величиной

$$\frac{1}{2} (t_n - t_{n-1}) (t_n - t_{n-2}) K.$$

Так как

$$\begin{aligned} |f'(x_{n-1})| &\geq |f'(x_0)| - |f'(x_0) - f'(x_{n-1})| \geq \frac{1}{B} - |x_{n-1} - x_0| K \geq \\ &\geq \frac{1}{B} - [|x_{n-1} - x_{n-2}| + \dots + |x_1 - x_0|] K \geq \\ &\geq \frac{1}{B} - K [(t_{n-1} - t_{n-2}) + \dots + (t_1 - t_0)] = \frac{1}{B} - K t_{n-1} \end{aligned}$$

для абсолютной величины знаменателя в (1.8.11) находим оценку снизу

$$\begin{aligned} |f'(x_{n-1}) + \frac{1}{2} (x_n - x_{n-1}) f''(\eta)| &\geq \frac{1}{B} - K t_{n-1} - \frac{1}{2} K (t_n - t_{n-1}) = \\ &= \frac{1}{B} - K \frac{t_n + t_{n-1}}{2}. \end{aligned}$$

Попутно отметим, что

$$K \frac{t_n + t_{n-1}}{2} - \frac{1}{B} = P' \left(\frac{t_n + t_{n-1}}{2} \right).$$

Так как $\frac{t_n + t_{n-1}}{2} < t^*$, то $P' \left(\frac{t_n + t_{n-1}}{2} \right) < 0$ и $\frac{1}{B} - K \frac{t_n + t_{n-1}}{2} > 0$.

Таким образом, из (1.8.11) получаем

$$|x_{n+1} - x_n| \leq \frac{\frac{1}{2} K (t_n - t_{n-1}) (t_n - t_{n-2})}{\frac{1}{B} - \frac{1}{2} K (t_n - t_{n-1})}. \quad (1.8.12)$$

Если проделать для уравнения $P(t) = 0$ вычисления, сходные с теми, которые были проделаны нами для получения (1.8.11), найдем:

$$\begin{aligned}
P(t_n) &= \frac{1}{2} K(t_n - t_{n-1})(t_n - t_{n-2}), \\
P(t_n) - P(t_{n-1}) &= \frac{1}{2} K(t_n^2 - t_{n-1}^2) - \frac{1}{B} (t_n - t_{n-1}) = \\
&= (t_n - t_{n-1}) \left[\frac{1}{B} - \frac{1}{2} K \left(\frac{t_n + t_{n-1}}{2} \right) \right], \\
t_{n+1} - t_n &= - \frac{(t_n - t_{n-1}) P(t_n)}{P(t_n) - P(t_{n-1})} = \frac{\frac{1}{2} K(t_n - t_{n-1})(t_n - t_{n-2})}{\frac{1}{B} - \frac{1}{2} K(t_n + t_{n-1})}.
\end{aligned}$$

Сравнение этого результата с (1.8.12) приводит к нужному неравенству $|x_{n+1} - x_n| \leq t_{n+1} - t_n$. Этим завершается индуктивное доказательство (1.8.10). Остановимся сначала на случае конечного N .

Предполагая возможным построение x_0, x_1, \dots, x_N , мы должны считать, что x_0, x_1, \dots, x_{N-1} лежат на отрезке $|x - x_0| \leq \delta$. Покажем, что x_N также принадлежит этому отрезку. Действительно,

$$\begin{aligned}
|x_N - x_0| &\leq |x_N - x_{N-1}| + |x_{N-1} - x_{N-2}| + \dots + |x_1 - x_0| \leq (t_N - t_{N-1}) + \\
&+ (t_{N-1} - t_{N-2}) + \dots + (t_1 - t_0) = t_N - t_0 = t_N < t^* = \frac{1 - \sqrt{1 - 2h}}{h} \eta \leq \delta.
\end{aligned}$$

Так как x_N принадлежит области определения $f(x)$, приближение x_{N+1} не может быть построено только по причине равенства $f(x_N) - f(x_{N-1}) = 0$. Оценим разность

$$|f(x_N) - f(x_{N-1})| = \left| (x_N - x_{N-1}) \left[f'(x_{N-1}) + \frac{1}{2} (x_N - x_{N-1}) f''(\eta) \right] \right|.$$

Применив к сумме, стоящей в прямых скобках, рассуждения, вполне аналогичные тем, которые были проделаны несколькими строками выше при оценке снизу абсолютной величины делителя в (1.8.11), получим неравенство

$$\begin{aligned}
|f(x_N) - f(x_{N-1})| &\geq |x_N - x_{N-1}| \left[\frac{1}{B} - K \frac{1}{2} (t_N + t_{N-1}) \right] = \\
&= |x_N - x_{N-1}| \cdot \left| P' \left[\frac{1}{2} (t_N + t_{N-1}) \right] \right|.
\end{aligned}$$

Производная $P' \left[\frac{1}{2}(t_{N-1} + t_N) \right]$ отлична от нуля ввиду $\frac{1}{2}(t_{N-1} + t_N) < t^*$.

Поэтому изучаемая разность может равняться нулю только в том случае, когда $x_N = x_{N-1}$. Но тогда, как отмечалось перед доказываемой теоремой, x_{N-1} есть решение заданного уравнения: $x_{N-1} = x^*$. Таким образом, в случае конечного N решение уравнения находится по методу секущих за конечное число шагов.*)

Перейдем к случаю $N = \infty$, когда последовательность x_n является бесконечной. Проверим ее сходимость:

$$\begin{aligned} |x_{n+p} - x_n| &= |(x_{n+p} - x_{n+p-1}) + (x_{n+p-1} - x_{n+p-2}) + \dots + (x_{n+1} - x_n)| \leq \\ &\leq (t_{n+p} - t_{n+p-1}) + (t_{n+p-1} - t_{n+p-2}) + \dots + (t_{n+1} - t_n) = t_{n+p} - t_n. \end{aligned}$$

Так как последовательность t_n сходится, для нее признак Больцано — Коши выполняется. Из полученного неравенства следует, что признак будет выполняться и для x_n и последовательность x_n также будет сходящейся: $x_n \rightarrow x^*$.

Неравенство (1.8.9), характеризующее скорость сходимости, сразу же получается, если в неравенстве $|x_{n+p} - x_n| \leq t_{n+p} - t_n$ перейти к пределу при $p \rightarrow \infty$ и заметить, что при этом $x_{n+p} \rightarrow x^*$ и $t_{n+p} \rightarrow t^*$.

Наконец, если в правиле (1.8.7) перейти к пределу при $n \rightarrow \infty$, то, ввиду $x_{n+1} \rightarrow x^*$ и $x_n \rightarrow x^*$, мы придем к заключению, что

$$\frac{(x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})} \rightarrow 0$$

и, так как

$$\frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} \rightarrow f'(x^*),$$

должно быть $f(x_n) \rightarrow 0$ независимо от значения $f'(x^*)$. Для предельного значения x^* , из-за непрерывности f в точке x^* , должно выполняться равенство $f(x^*) = 0$ и x^* должно быть решением уравнения $f(x) = 0$. Теорема 1 доказана.

2. Видоизменение с начальным значением производной.

В этом последнем видоизменении метода Ньютона, как и в предыдущем, освобождаются от вычисления значений производных $f'(x_n)$ на каждом шаге. Пользуются только одним начальным значением $f'(x_0)$. Последовательные приближения вычисляются по правилу

) Аналогичное может случиться и при вычислениях по основному методу Ньютона, с тем лишь различием, что при $x_{N-1} = x^$ и $f'(x_{N-1}) \neq 0$ расчетная формула не теряет смысла, вычисления могут быть продолжены, но последовательность приближений становится стационарной: $x_{N-1} = x_N = \dots = x^*$. По этой причине в § 1.7 мы не останавливались на этом исключительном случае.

$$x'_{n+1} = x'_n - \frac{f(x'_n)}{f'(x'_0)} \quad (n=0, 1, 2, \dots). \quad (1.8.13)$$

Исходное значение x'_0 считается известным. Геометрический смысл правила состоит в следующем: на графике l функции $f(x)$ возьмем точку $M_n[x'_n, f(x'_n)]$ и через нее проведем прямую линию, имеющую тот же тангенс угла наклона к оси Ox , что и касательная к l в начальной точке $M_0[x'_0, f(x'_0)]$. Уравнение такой прямой есть

$$y - f(x'_n) = f'(x'_0)(x - x'_n).$$

Правило (1.8.13) означает, что за следующее приближение x'_{n+1} принимается абсцисса точки пересечения указанной прямой с осью Ox .

Будем считать x'_n и x'_{n+1} близкими к решению x^* и рассмотрим погрешности $\varepsilon'_n = x^* - x'_n$, $\varepsilon'_{n+1} = x^* - x'_{n+1}$. Правило (1.8.13) дает уравнение для погрешностей

$$\varepsilon'_{n+1} = \varepsilon'_n + \frac{f(x^* - \varepsilon'_n)}{f'(x'_0)}.$$

Если воспользоваться разложением

$$f(x^* - \varepsilon'_n) = -\varepsilon'_n f'(x^*) + \frac{1}{2} \varepsilon'^2_n f''(x^*), \dots$$

подставить его в правую часть уравнения и сохранить только линейные члены, получим приближенное равенство, дающее описание изменения погрешности ε'_n вблизи точного решения:

$$\varepsilon'_{n+1} \approx \varepsilon'_n \left[1 - \frac{f'(x^*)}{f'(x'_0)} \right]. \quad (1.8.14)$$

Закон изменения близок к геометрической прогрессии со знаменателем $q = 1 - \frac{f'(x^*)}{f'(x'_0)}$. Так как исходное приближение x'_0 обычно берется близким к решению x^* , отношение $\frac{f'(x^*)}{f'(x'_0)}$ является близким к единице и знаменатель q , как правило, имеет небольшое значение.

Сравнение (1.8.14) с законом (1.7.4) убывания погрешностей в основном методе Ньютона позволяет сказать, что видоизменение метода с начальным значением производной имеет сходимость более медленную, чем в основном методе.

Теорема 2. Пусть для уравнения $f(x) = 0$ и начального значения x'_0 выполняются условия:

1) $f(x)$ определена на отрезке $|x - x_0'| \leq \delta$ и дважды непрерывно дифференцируема там;

$$2) f'(x_0') \neq 0 \text{ и } \frac{1}{|f'(x_0')|} \leq B;$$

3) выполнено условие

$$\left| \frac{f(x_0')}{f'(x_0')} \right| \leq \eta;$$

4) вторая производная $f''(x)$ ограничена по абсолютному значению на $|x - x_0'| \leq \delta$ числом K :

$$|f''(x)| \leq K;$$

5) для чисел B, η, K выполнено условие

$$h = BK\eta \leq \frac{1}{2};$$

6) верно неравенство

$$\frac{1 - \sqrt{1 - 2h}}{h} \eta \leq \delta.$$

Тогда:

1) приближение x_n' , определяемое правилом (1.8.13), может быть построено для любого n ;

2) последовательность приближений x_n' сходится:

$$x_n' \rightarrow x^* \quad (n \rightarrow \infty);$$

3) предельное значение x^* есть решение заданного уравнения;

4) скорость сходимости оценивается неравенством

$$|x^* - x_n'| \leq t^* - t_n',$$

где t^* есть меньший корень квадратного уравнения

$$P(t) = \frac{1}{2} K t^2 - \frac{1}{B} t + \frac{\eta}{B} = 0$$

и t_n' есть последовательность приближений к нему, построенная по пра-

вилу $t'_{n+1} = t_n' - \frac{P(t_n')}{P'(t_n')}$ при $t_0' = 0$.

Доказательство. Приближение t_n' будет монотонно возрастать при увеличении n и стремиться к t^* . Убедиться в этом можно при помощи чертежа.

Покажем, что приближения x_n' могут быть построены при всяком n , принадлежат области $|x - x_0'| \leq \delta$ и для x_n' выполняется неравенство

$$|x'_{n+1} - x_n'| \leq t'_{n+1} - t_n' \quad (n=0, 1, 2, \dots). \quad (1.8.15)$$

При $n=0$ неравенство верно, так как $x_1' = x_0' - \frac{f(x_0')}{f'(x_0')}$ может быть построено и $|x_1' - x_0'| = \left| \frac{f(x_0')}{f'(x_0')} \right| \leq \eta$. Кроме того, $t_1' = t_0' - \frac{P(t_0')}{P'(t_0')}$, $t_1' - t_0' = t_1' = \eta$ и, стало быть, $|x_1' - x_0'| \leq t_1' - t_0'$.

Наконец, $|x_1' - x_0'| \leq \eta \leq \frac{1 - \sqrt{1-2h}}{h} \eta \leq \delta$ и x_1' принадлежит области $|x - x_0'| \leq \delta$.

Допустим, что x_0', x_1', \dots, x_n' могут быть построены, принадлежат области $|x' - x_0'| \leq \delta$ и для них выполнены неравенства $|x'_{k+1} - x_k'| \leq t'_{k+1} - t_k' \quad (k=0, 1, \dots, n-1)$. Приближение $x'_{n+1} = x_n' - \frac{f(x_n')}{f'(x_0')}$, очевидно, может быть построено, так как x_n' принадлежит области определения f и $f'(x_0') \neq 0$. Займемся оценкой разности $x'_{n+1} - x_n'$. На основании правила (1.8.13) и при помощи простых преобразований получаем для нее приводимое ниже представление:

$$\begin{aligned} x'_{n+1} - x_n' &= x_n' - x'_{n-1} - \frac{f(x_n') - f(x'_{n-1})}{f'(x_0')} = \\ &= -\frac{1}{f'(x_0')} [f(x_n') - f(x'_{n-1}) - f'(x_0')(x_n' - x'_{n-1})] = \\ &= -\frac{1}{f'(x_0')} [f(x_n') - f(x'_{n-1}) - f'(x'_{n-1})(x_n' - x'_{n-1})] - \\ &\quad - \frac{1}{f'(x_0')} (x_n' - x'_{n-1}) [f'(x'_{n-1}) - f'(x_0')]. \end{aligned}$$

В первой квадратной скобке последней части равенства стоит остаток тейлорова разложения около точки x'_{n-1} . Для него известно интегральное представление через вторую производную f'' :

$$\int_{x'_{n-1}}^{x_n'} f''(t) (x_n' - t) dt,$$

которое дает возможность оценить первую из скобок величиной

$$\frac{1}{2} K |x_n' - x_{n-1}'|^2.$$

Для второй квадратной скобки

$$|f'(x_{n-1}') - f'(x_0')| = \left| \int_{x_0'}^{x_{n-1}'} f''(t) dt \right| \leq K |x_{n-1}' - x_0'|.$$

На основании индуктивного предположения

$$|x_n' - x_{n-1}'| \leq t_n' - t_{n-1}'$$

и

$$|x_{n-1}' - x_0'| \leq |x_{n-1}' - x_{n-2}'| + \dots + |x_1' - x_0'| \leq (t_{n-1}' - t_{n-2}') + \dots + (t_1' - t_0') = t_{n-1}' - t_0' = t_{n-1}'.$$

Наконец, так как $\left| \frac{1}{f'(x_0')} \right| \leq B$, найдем

$$\begin{aligned} |x_{n+1}' - x_n'| &\leq BK \left[\frac{1}{2} (t_n' - t_{n-1}')^2 + (t_n' - t_{n-1}') t_{n-1}' \right] = \\ &= \frac{1}{2} BK (t_n'^2 - t_{n-1}'^2). \end{aligned}$$

Сходные вычисления для уравнения $P(t) = 0$ дадут

$$t_{n+1}' - t_n' = \frac{1}{2} BK (t_n'^2 - t_{n-1}'^2),$$

и из сличения этого равенства с предыдущей оценкой получается нужное неравенство

$$|x_{n+1}' - x_n'| \leq t_{n+1}' - t_n'.$$

Проверим принадлежность x_{n+1} области $|x - x_0'| \leq \delta$:

$$\begin{aligned} |x_{n+1}' - x_0'| &\leq |x_{n+1}' - x_n'| + \dots + |x_1' - x_0'| \leq (t_{n+1}' - t_n') + \\ &+ \dots + (t_1' - t_0') = t_{n+1}' - t_0' < t^* \leq \delta. \end{aligned}$$

Дальнейшие рассуждения будут весьма сходными с концом доказательства теоремы 1 и мы приведем их очень коротко. При помощи (1.8.15) легко получается неравенство $|x_{n+p}' - x_n'| \leq t_{n+p}' - t_n'$. Так как последовательность t_n' сходится, то для нее выполняется признак Больцано — Коши. Из неравенства следует, что признак будет выполняться и для

последовательности x_n' и она будет сходящейся: $\lim x_n' = x^*$. Если же в неравенстве перейти к пределу при $p \rightarrow \infty$, мы получим утверждение (4) теоремы об оценке скорости сходимости. Наконец, если в правиле вычислений (1.8.13) перейти к пределу при $n \rightarrow \infty$, мы убедимся в том, что $\lim f(x_n') = 0$. Ввиду $x_n' \rightarrow x^*$ и непрерывности f это дает $f(x^*) = 0$ и x^* есть, следовательно, решение заданного уравнения.

Оценка скорости сходимости, указанная в утверждении (4), является неулучшаемой в условиях теоремы, так как она достигается для квадратного уравнения $P(t) = 0$. Можно указать другую, более простую и наглядную, но менее точную оценку. Приближения t_n и t_{n-1} связаны зависимостью

$$t_n' = t'_{n-1} - \frac{P(t'_{n-1})}{P'(0)} = \eta + \frac{1}{2} BK t_{n-1}^2.$$

Для точного решения t^* справедливо равенство

$$t^* = \eta + \frac{1}{2} BK t^{*2}.$$

Вычитая эти равенства почленно, найдем

$$t^* - t_n' = \frac{1}{2} BK (t^{*2} - t_{n-1}^2) = \frac{1}{2} BK (t^* + t'_{n-1}) (t^* - t'_{n-1}),$$

и так как $t'_{n-1} < t^* = \frac{1 - \sqrt{1-2h}}{h} \eta$,

$$t^* - t_n' < BK t^* (t^* - t'_{n-1}) = (1 - \sqrt{1-2h}) (t^* - t'_{n-1}).$$

Применим это неравенство n раз:

$$t^* - t_n' < q^n t^*, \quad q = 1 - \sqrt{1-2h} < 1 \quad \text{при} \quad h < \frac{1}{2}.$$

Последнее неравенство говорит о том, что сходимость $t_n' \rightarrow t^*$ происходит не медленнее, чем со скоростью геометрической прогрессии со знаменателем q .

§ 1.9. ОПЕРАТОРНЫЕ УРАВНЕНИЯ И МЕТОД НЬЮТОНА

В § 1.7 и 1.8 мы рассматривали теорию метода Ньютона для уравнения с одной численной неизвестной величиной. Основная же идея этого метода о приведении нелинейного уравнения к последовательности линейных уравнений имеет, очевидно, много более общее значение: она применима к системам уравнений с несколькими численными неизвестными, к уравнениям, в которых роль неизвестной величины играет не число, а функция, например, к нелинейным дифференциальным и интегральным уравнениям, и многим другим.

Было бы желательно построить теорию метода Ньютона, не зависящую от конкретного вида уравнения и опирающуюся только на абстрактные математические понятия. Это оказалось возможным сделать в достаточно общем виде, если воспользоваться некоторыми понятиями и результатами функционального анализа.*)

Докажем сейчас две теоремы, достаточно общие и удобные для многих приложений, об осуществимости основного метода Ньютона для операторных уравнений, сходимости последовательных приближений к точному решению и о единственности решения.**)

Теоремы сформулированы в форме, сходной с аналогичными теоремами § 1.7 для численного уравнения. Доказательства теорем в значительной мере повторяют рассуждения, проведенные в указанном параграфе.

Пусть X и Y есть два полных, линейных, нормированных пространства, элементы которых будем обозначать x и y соответственно.

Предположим, что в некоторой области D пространства X определен нелинейный оператор $y=f(x)$, значения которого принадлежат пространству Y . Оператор $f(x)$ предполагается дважды дифференцируемым в D в смысле Фреше (добавление I, § 3);

Рассмотрим уравнение

$$f(x)=0. \quad (1.9.1)$$

В правой части его стоит нулевой элемент из Y . Будем считать, что мы знаем исходное приближение x_0 к решению уравнения (1.9.1). Правила вычислений, позволяющие найти следующие приближения по предыдущим, строятся на том же основании, как и в случае одного численного уравнения. Полагая $f(x)=f(x_0)+[f(x)-f(x_0)]$, возьмем вместо разности $f(x)-f(x_0)$ дифференциал оператора на элементе x_0 : $f'(x_0)(x-x_0)$ и заменим заданное уравнение $f(x)=0$ приближенно линейным уравнением

$$f(x) \approx f(x_0) + f'(x_0)(x-x_0) = 0. \quad (1.9.2)$$

Решая его, найдем улучшенное приближенное значение x_1 . Если для $f'(x_0)$ существует обратный оператор $[f'(x_0)]^{-1}$, переводящий Y в X , для x_1 можно написать следующее явное выражение:

$$x_1 = x_0 - [f'(x_0)]^{-1} f(x_0).$$

Повторяя для x_1 те же операции, построим по x_1 второе улучшенное приближение x_2 и т. д. В общем виде правило нахождения следующего приближения по предыдущему будет

$$x_{n+1} = x_n - [f'(x_n)]^{-1} f(x_n) \quad (n=0, 1, 2, \dots). \quad (1.9.3)$$

Условием возможности построения последовательности x_n является выполнение двух требований:

- 1) принадлежность x_n ($n=0, 1, \dots$) к области D определения оператора f ;
- 2) существование обратных операторов $[f'(x_n)]^{-1}$.

Докажем теорему о разрешимости уравнения $f(x)=0$, об осуществимости процесса Ньютона (1.9.3) и о сходимости последовательности x_n к решению уравнения.

Теорема 1. Пусть выполнены условия:

- 1) оператор $f(x)$ определен в замкнутом шаре

$$\|x-x_0\| \leq \delta \quad (\alpha)$$

около исходного приближения x_0 , дважды дифференцируем там и вторая производная от него по норме ограничена в этом шаре числом K :

*) Необходимые сведения из функционального анализа можно найти в добавлении I.

**) Более полное ознакомление с общей теорией метода Ньютона можно найти в книгах [3 и 4].

$$\|f''(x)\| \leq K, \quad \|x - x_0\| \leq \delta;$$

2) оператор $f'(x_0)$ имеет обратный $\Gamma_0 = [f'(x_0)]^{-1}$ и известна оценка его нормы:

$$\|\Gamma_0\| \leq B;$$

3) на начальном элементе x_0 соблюдено неравенство

$$\|\Gamma_0 f(x_0)\| \leq \eta;$$

4) для B, K, η выполнено условие

$$h = BK\eta \leq \frac{1}{2};$$

5) для δ верно неравенство

$$\frac{1 - \sqrt{1 - 2h}}{h} \eta \leq \delta.$$

Тогда:

1) уравнение $f(x) = 0$ имеет в области (α) решение;

2) последовательные приближения x_n ($n=0, 1, 2, \dots$) процесса Ньютона могут быть построены, принадлежат области (α) и сходятся к решению уравнения: $\lim_{n \rightarrow \infty} x_n = x^*$;

3) быстрота сходимости оценивается неравенством

$$\|x^* - x_n\| \leq t^* - t_n, \quad (1.9.4)$$

где t_n и t^* имеют тот же смысл, что и в теореме 1, § 1.7.

Доказательство. Покажем, что при условиях теоремы приближения x_n могут быть построены при любых значениях n , все x_n принадлежат области (α) и для них выполняется неравенство

$$\|x_{n+1} - x_n\| \leq t_{n+1} - t_n. \quad (1.9.5)$$

Прибегнем к индукции. Так как x_0 принадлежит области (α) и по предположению (2) оператор $\Gamma_0 = [f'(x_0)]^{-1}$ существует, первое улучшенное приближение $x_1 = x_0 - \Gamma_0 f(x_0)$ может быть построено. Кроме того, по предположению (3) $\|x_1 - x_0\| = \|\Gamma_0 f(x_0)\| \leq \eta$. Если же обратить внимание на то, что $t_1 - t_0 = -\frac{p(t_0)}{p'(t_0)} = -\frac{p(0)}{p'(0)} = \eta$, то мы убедимся, что неравенство (1.9.5) для x_0 и x_1 выполняется.

Допустим теперь, что приближения x_0, x_1, \dots, x_n могут быть построены, принадлежат (α) и для них выполняются неравенства $\|x_{k+1} - x_k\| \leq t_{k+1} - t_k$ ($k=0, 1, \dots, n-1$). Мы покажем, что приближение x_{n+1} может быть построено по правилу (1.9.3), если установим, что существует оператор $[f'(x_n)]^{-1}$. Для этого воспользуемся следующей простой леммой об обратном операторе.

Лемма. Пусть линейный оператор H преобразует полное нормированное линейное пространство X в себя. Если $\|H\| = q < 1$, то оператор $I - H$ имеет обратный $(I - H)^{-1}$,

при этом $\|(I - H)^{-1}\| \leq \frac{1}{1 - q}$.

Доказательство. Ряд $Ax = x + Hx + H^2x + \dots$ сходится, так как $\|H^n x\| \leq q^n \|x\|$ и $\|Ax\| \leq \frac{1}{1 - q} \|x\|$. Поэтому $\|A\| \leq \frac{1}{1 - q}$. Кроме того,

$$(I-H)Ax = (x+Hx+H^2x+\dots) - (Hx+H^2x+\dots) = x,$$

$$A(I-H)x = (x-Hx) + (Hx-H^2x) + \dots = x.$$

Отсюда следует $A = (I-H)^{-1}$ и $\|(I-H)^{-1}\| \leq \frac{1}{1-q}$, что доказывает лемму.

Рассмотрим теперь оператор $H = \Gamma_0[f'(x_n) - f'(x_0)]$, дающий отображение X в себя, и оценим его норму. В теории операторов известна приводимая ниже теорема об оценке изменения оператора.*)

Если $F(x)$ — дифференцируемый оператор, то для его изменения верно неравенство

$$\|F(x+\Delta x) - F(x)\| \leq \sup_{0 \leq \theta \leq 1} \|F'(x+\theta\Delta x)\| \|\Delta x\|.$$

Применим эту оценку к $f'(x_n) - f'(x_0)$:

$$\|f'(x_n) - f'(x_0)\| \leq \sup_{0 \leq \theta \leq 1} \|f''[x_0 + \theta(x_n - x_0)]\| \|x_n - x_0\| \leq K \|x_n - x_0\|.$$

Но

$$\|x_n - x_0\| \leq \|x_n - x_{n-1}\| + \|x_{n-1} - x_{n-2}\| + \dots + \|x_1 - x_0\|,$$

и так как по индуктивному допущению

$$\|x_{k+1} - x_k\| \leq t_{k+1} - t_k \quad (k=0, 1, \dots, n-1),$$

то

$$\|x_n - x_0\| \leq (t_n - t_{n-1}) + (t_{n-1} - t_{n-2}) + \dots + (t_1 - t_0) = t_n$$

и

$$\|f'(x_n) - f'(x_0)\| \leq K t_n.$$

Поэтому

$$\|H\| \leq \|\Gamma_0\| \|f'(x_n) - f'(x_0)\| \leq K B t_n < K B t_n^* = K B \eta \frac{1 - \sqrt{1-2h}}{h} = 1 - \sqrt{1-2h} \leq 1$$

и

$$\|H\| \leq K B t_n < 1.$$

Ввиду приведенной выше леммы оператор

$$I+H = I + \Gamma_0[f'(x_n) - f'(x_0)]$$

имеет обратный

$$(I+H)^{-1} = \{I + \Gamma_0[f'(x_n) - f'(x_0)]\}^{-1}$$

и верна оценка нормы

$$\|(I+H)^{-1}\| \leq \frac{1}{1 - B K t_n^*}.$$

Положим $\Gamma_n = (I+H)^{-1} \Gamma_0$ и, воспользовавшись известным равенством $(AB)^{-1} = B^{-1}A^{-1}$, получим

$$\begin{aligned} \Gamma_n &= (I+H)^{-1} \Gamma_0 = \{\Gamma_0^{-1}[I + \Gamma_0(f'(x_n) - f'(x_0))]\}^{-1} = \\ &= \{f'(x_0) + [f'(x_n) - f'(x_0)]\}^{-1} = \{f'(x_n)\}^{-1}. \end{aligned}$$

*) Доказательство ее указано в добавлении I.

Мы доказали существование оператора $\{f'(x_n)\}^{-1} = \Gamma_n$ и этим установили возможность построения x_{n+1} . Легко получается оценка

$$\|\{f'(x_n)\}^{-1}\| = \|\Gamma_n\| \leq \| (I+H)^{-1} \| \quad \|\Gamma_0\| \leq \frac{B}{1-BKt_n} = \frac{1}{\frac{1}{B} - Kt_n} = -\frac{1}{P'(t_n)}.$$

Под $P(t)$ здесь понимается квадратный многочлен $P(t) = \frac{1}{2} Kt^2 - \frac{1}{B}t + \frac{\eta}{B}$, который участвовал в теореме 1, §1.7.

Нам осталось еще найти оценку $f(x_n)$. По предположению, x_n находится по правилу (1.9.3) и для него верно равенство

$$x_n = x_{n-1} - [f'(x_{n-1})]^{-1} f(x_{n-1})$$

или

$$f(x_{n-1}) + f'(x_{n-1})(x_n - x_{n-1}) = 0.$$

Поэтому

$$f(x_n) = f(x_n) - f(x_{n-1}) - f'(x_{n-1})(x_n - x_{n-1}).$$

Это выражение можно рассматривать как разницу между изменением $f(x_n) - f(x_{n-1})$ оператора f и его дифференциалом $f'(x_{n-1})(x_n - x_{n-1})$. В функциональном анализе известна теорема об оценке такой разности.

Если оператор F двукратно дифференцируем, то для него верно неравенство *)

$$\|F(x+\Delta x) - F(x) - F'(x)\Delta x\| \leq \frac{1}{2} \max_{0 \leq \theta \leq 1} \|F''(x+\theta\Delta x)\| \|\Delta x\|^2.$$

Применим это неравенство к приведенному выше выражению $f(x_n)$ и воспользуемся тем, что $\|x_n - x_{n-1}\| \leq t_n - t_{n-1}$:

$$\|f(x_n)\| \leq \frac{1}{2} \max_{0 \leq \theta \leq 1} \|f''[x_{n-1} + \theta(x_n - x_{n-1})]\| \|x_n - x_{n-1}\|^2 \leq \frac{1}{2} K(t_n - t_{n-1})^2.$$

При доказательстве теоремы 1 §1.7 обращалось внимание на то, что

$$\frac{1}{2} K(t_n - t_{n-1})^2 = P(t_n).$$

Следовательно,

$$\|x_{n+1} - x_n\| \leq \|\{f'(x_n)\}^{-1}\| \|f(x_n)\| \leq -\frac{P(t_n)}{P'(t_n)} = t_{n+1} - t_n$$

и неравенство (1.9.5) для x_n и x_{n+1} действительно выполняется.

Проверка принадлежности x_{n+1} внутренности шара не вызывает затруднений:

$$\begin{aligned} \|x_{n+1} - x_0\| &\leq \|x_{n+1} - x_n\| + \|x_n - x_{n-1}\| + \dots + \|x_1 - x_0\| \leq (t_{n+1} - t_n) + \\ &+ (t_n - t_{n-1}) + \dots + (t_1 - t_0) = t_{n+1} - t_0 = t_{n+1} < t^* = \frac{1 - \sqrt{1-2h}}{h} \eta \leq \delta. \end{aligned}$$

*) Доказательство приведено в добавлении I.

После того как установлено неравенство (1.9.5), дальнейший ход доказательства теоремы протекает без труда. Так как последовательность t_n сходится, для нее признак Больцано — Коши выполняется, а из оценки

$$\begin{aligned} \|x_{n+p} - x_n\| &\leq \|x_{n+p} - x_{n+p-1}\| + \|x_{n+p-1} - x_{n+p-2}\| + \dots + \|x_{n+1} - x_n\| \leq \\ &\leq (t_{n+p} - t_{n+p-1}) + (t_{n+p-1} - t_{n+p-2}) + \dots + (t_{n+1} - t_n) = t_{n+p} - t_n \end{aligned}$$

следует, что признак выполняется и для последовательности x_n . Ввиду же полноты пространства X , последовательность x_n будет сходиться: $\lim x_n = x^*$. Все x_n принадлежат замкнутому шару (α) и предел x^* также будет ему принадлежать.

Для доказательства того, что предел x^* есть решение уравнения $f(x) = 0$, рассмотрим правило (1.9.3), взяв его в форме $f(x_n) = f'(x_n)(x_{n+1} - x_n)$ и перейдем к пределу при $n \rightarrow \infty$. При этом будет $x_{n+1} - x_n \rightarrow 0$, и так как f и f' дифференцируемы и, следовательно, непрерывны всюду в шаре (α) , в частности на элементе x^* , то $f(x_n)$ и $f'(x_n)$ будут стремиться соответственно к $f(x^*)$ и $f'(x^*)$. В пределе будем иметь $f(x^*) = 0$.

Утверждение об оценке скорости сходимости (1.9.4) получится из неравенства $\|x_{n+p} - x_n\| \leq t_{n+p} - t_n$, если в нем фиксировать n и устремить p к $+\infty$.

Отметим, что оценка (1.9.5) в условиях теоремы не может быть улучшена, так как она достигается для численного уравнения $P(t) = 0$.

В конце § 1.7 было показано, что

$$t^* - t_n \leq \frac{1}{2^{n-1}} (2h)^{2^{n-1}},$$

и оценка (1.9.4) может быть поэтому заменена более наглядной, но несколько более грубой оценкой

$$\|x^* - x_n\| \leq \frac{1}{2^{n-1}} (2h)^{2^{n-1}} \eta. \quad (1.9.6)$$

Перейдем теперь к проблеме единственности решения.

Теорема 2. Пусть для оператора $f(x)$ и начального элемента x_0 выполнены условия

$$(1)-(5) \text{ теоремы 1 с тем различием, что } \delta \geq \frac{1 + \sqrt{1-2h}}{h} \eta = t^{**}.$$

Тогда уравнение $f(x) = 0$ имеет единственное решение при $h < \frac{1}{2}$ в области $\|x - x_0\| < t^{**}$ и при $h = \frac{1}{2}$ в области $\|x - x_0\| \leq t^{**}$.

Доказательство. Ввиду того что условия теоремы 1 выполняются, уравнение имеет решение x^* в области $\|x - x_0\| \leq t^* \leq t^{**}$, составляющей часть области, предусмотренной теоремой 2. Нам нужно показать, что всякое другое решение будет совпадать с x^* .

Рассмотрим случай $h < \frac{1}{2}$ и допустим, что существует решение \bar{x} в области $\|x - x_0\| < t^{**}$. Можно положить

$$\|\bar{x} - x_0\| \leq \theta t^{**} = \theta(t^{**} - t_0) \quad (0 \leq \theta < 1). \quad (1.9.7)$$

Для сокращения записи введем оператор $F(x) = x - \Gamma_0 f(x)$ и отметим несколько его свойств, которыми мы воспользуемся в преобразованиях:

$$F(\bar{x}) = \bar{x}, \quad F'(x_0) = 0, \quad F''(x) = -\Gamma_0 f''(x).$$

Выполним оценки

$$\begin{aligned}\|\bar{x}-x_1\| &= \|\bar{x}-[x_0-\Gamma_0 f(x_0)]\| = \|F(\bar{x})-F(x_0)\| = \|F(\bar{x})-F(x_0)-F'(x_0)(\bar{x}-x_0)\| = \\ &= \frac{1}{2} \sup_{0 \leq \theta \leq 1} \|F''[x_0+\theta(\bar{x}-x_0)]\| \|\bar{x}-x_0\|^2 = \frac{1}{2} \sup_{0 \leq \theta \leq 1} \|\Gamma_0 f''[x_0+\theta(\bar{x}-x_0)]\| \|\bar{x}-x_0\|^2 \leq \\ &\leq \frac{1}{2} BK \|\bar{x}-x_0\|^2 \leq \frac{1}{2} BK \theta^2 t^{**2}\end{aligned}$$

и так как

$$\frac{1}{2} BK t^{**2} = -\frac{1}{P'(t_0)} [P(t^{**}) - P(t_0) - (t^{**}-t_0)P'(t_0)] = t^{**} - \left(t_0 - \frac{P(t_0)}{P'(t_0)}\right) = t^{**} - t_1,$$

то

$$\|\bar{x}-x_1\| \leq \theta^2 (t^{**}-t_1).$$

Сравнение последней оценки с (1.9.7) говорит о том, что она получается из (1.9.7) заменой θ на θ^2 и t_0 на t_1 . n -кратное применение этого правила приведет к неравенству

$$\|\bar{x}-x_n\| \leq \theta^{2n} (t^{**}-t_n) < \theta^{2n} t^{**}. \quad (1.9.8)$$

$\theta < 1$ и не зависит от n , поэтому $\|\bar{x}-x_n\| \rightarrow 0$ и $x_n \rightarrow \bar{x}$. Но по теореме 1 $x_n \rightarrow x^*$ и, следовательно, $\bar{x} = x^*$.

При $h = \frac{1}{2}$ число θ может равняться единице. Но тогда $t^{**} = t^*$ и ввиду $t_n \rightarrow t^*$ из (1.9.8) вытекает $\|\bar{x}-x_n\| \rightarrow 0$, откуда, как выше, следует $\bar{x} = x^*$.

§ 1.10. МЕТОД НЬЮТОНА ДЛЯ СИСТЕМ УРАВНЕНИЙ

Рассмотрим систему v уравнений с v численными неизвестными x_1, x_2, \dots, x_v :

$$f_1(x_1, x_2, \dots, x_v) = 0, \quad f_2(x_1, x_2, \dots, x_v) = 0, \dots, f_v(x_1, x_2, \dots, x_v) = 0. \quad (1.10.1)$$

Для сокращения записи введем v -мерное векторное пространство, элементами которого будут упорядоченные совокупности v чисел $x = (x_1, x_2, \dots, x_v)$. Функцию $y_i = f_i(x_1, x_2, \dots, x_v)$ ниже часто будем записывать $y_i = f_i(x)$.

Одновременно введем v -мерную вектор-функцию $f(x) = (f_1(x), f_2(x), \dots, f_v(x))$. Система (1.10.1) коротко запишется

$$f(x) = 0. \quad (1.10.2)$$

Метод Ньютона для системы (1.10.1) является естественным обобщением этого метода для одного численного уравнения, который был рассмотрен в § 1.7. Как и там, основная идея метода будет состоять в приведении нелинейной системы (1.10.1) к решению последовательности линейных систем. Такое приведение достигается путем выделения из уравнений системы линейных частей, являющихся главными при малых погрешностях.

Пусть нам известно приближение $x^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_v^{(0)})$ к решению системы. Для наших целей удобнее рассматривать не точное вектор-решение $x = (x_1, \dots, x_v)$ системы, а вектор-погрешность

$$x - x^{(0)} = (x_1 - x_1^{(0)}, \dots, x_v - x_v^{(0)}) = \varepsilon = (\varepsilon_1, \dots, \varepsilon_v).$$

Уравнения для определения ε_i будут

$$f_i(x^{(0)} + \varepsilon) = 0 \quad (i=1, 2, \dots, v).$$

Разложив левую часть по степеням $\varepsilon_1, \dots, \varepsilon_v$ при помощи формулы Тейлора и сохранив лишь линейные члены, получим приближенную линейную систему

$$\sum_{j=1}^v \frac{\partial f_i(x^{(0)})}{\partial x_j} \varepsilon_j \approx -f_i(x^{(0)}) \quad (i=1, 2, \dots, v); \quad (1.10.3)$$

из которой мы сможем найти не точные, а лишь приближенные значения погрешностей, которые мы назовем $(\varepsilon_1^{(0)}, \dots, \varepsilon_v^{(0)}) = \bar{\varepsilon}^{(0)}$. Прибавляя $\bar{\varepsilon}_i^{(0)}$ к исходным значениям неизвестных, получим улучшенные их значения:

$$x_1^{(1)} = x_1^{(0)} + \varepsilon_1^{(0)}, \dots, x_v^{(1)} = x_v^{(0)} + \varepsilon_v^{(0)}. \quad (1.10.4)$$

Вектор $x^{(1)}$ мы в свою очередь можем улучшить, составив для него систему вида (1.10.3).

Из нее мы найдем главные части погрешности $\bar{\varepsilon}^{(1)} = x - x^{(1)}$ и т. д.

Каждое следующее приближение $x^{(n+1)}$ к решению будет находиться из линейной системы, составляемой по предшествующему приближению $x^{(n)}$,

$$\sum_{j=1}^v \frac{\partial f_i(x^{(n)})}{\partial x_j} (x_j^{(n+1)} - x_j^{(n)}) = -f_i(x^{(n)}) \quad (i=1, 2, \dots, v, \quad n=0, 1, 2, \dots). \quad (1.10.5)$$

Матрицей системы является значение матрицы Якоби

$$f'(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_v} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_v} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f_v}{\partial x_1} & \frac{\partial f_v}{\partial x_2} & \dots & \frac{\partial f_v}{\partial x_v} \end{bmatrix} \quad (1.10.6)$$

при $x = x^{(n)}$.

Система будет разрешимой и определенной только в том случае, когда определитель ее отличен от нуля: $D[f'(x^{(n)})] \neq 0$. Будем говорить, что процесс Ньютона для системы (1.10.1) осуществим, если система (1.10.5) может быть составлена и будет однозначно разрешимой при всех $n=0, 1, 2, \dots$. Условием осуществимости процесса, если считать функции f_i дифференцируемыми, является выполнение требований:

- 1) принадлежность $x^{(n)}$ ($n=0, 1, \dots$) области определения всех f_i ;
- 2) неравенство нулю определителей Якоби $D[f'(x^{(n)})]$ при $n=0, 1, 2, \dots$.

Изучение процесса Ньютона начнем с выяснения наглядной картины изменения погрешности $x^* - x^{(n)} = \varepsilon^{(n)} = (\varepsilon_1^{(n)}, \dots, \varepsilon_v^{(n)})$ вблизи точного решения, которое мы обозначим $x^* = (x_1^*, \dots, x_v^*)$. Погрешности $\varepsilon_1^{(n)}, \dots, \varepsilon_v^{(n)}$ будут малыми величинами. Подставим в (1.10.5) вместо $x_j^{(n)}$ и $x_j^{(n+1)}$ их выражения через погрешности $x_j^{(n)} = x_j^* - \varepsilon_j^{(n)}$ и $x_j^{(n+1)} = x_j^* - \varepsilon_j^{(n+1)}$. Функции f_i будем считать дважды непрерывно дифференцируемыми. Если принять во внимание равенства

$$x_j^{(n+1)} - x_j^{(n)} = \varepsilon_j^{(n)} - \varepsilon_j^{(n+1)}, \quad f_i(x^*) = 0,$$

$$f_i(x^{(n)}) = f_i(x^* - \varepsilon^{(n)}) = - \sum_{j=1}^v \frac{\partial f_i(x^*)}{\partial x_j} \varepsilon_j^{(n)} + \frac{1}{2} \sum_{j,k=1}^v \frac{\partial^2 f_i(x^*)}{\partial x_j \partial x_k} \varepsilon_j^{(n)} \varepsilon_k^{(n)} + \dots,$$

$$\frac{\partial f_i(x^{(n)})}{\partial x_j} = \frac{\partial f_i(x^* - \varepsilon^{(n)})}{\partial x_j} = \frac{\partial f_i(x^*)}{\partial x_j} - \sum_{k=1}^v \frac{\partial^2 f_i(x^*)}{\partial x_j \partial x_k} \varepsilon_k^{(n)} + \dots$$

и сохранить в результате подстановки лишь главные члены, то, пользуясь предположением о малости $\varepsilon_j^{(n)}$, получим

$$\sum_{j=1}^v \frac{\partial f_i(x^*)}{\partial x_j} \varepsilon_j^{(n+1)} \approx - \frac{1}{2} \sum_{j,k=1}^v \frac{\partial^2 f_i(x^*)}{\partial x_j \partial x_k} \varepsilon_j^{(n)} \varepsilon_k^{(n)} \quad (i=1, 2, \dots, v).$$

Можно воспользоваться матрицей $f'(x^*)$ и систему записать более просто:

$$f'(x^*) \varepsilon^{(n+1)} = - \frac{1}{2} \sum_{j,k=1}^v \frac{\partial^2 f(x^*)}{\partial x_j \partial x_k} \varepsilon_j^{(n)} \varepsilon_k^{(n)}. \quad (1.10.7)$$

Когда $D[f'(x^*)] \neq 0$ и матрица Якоби $f'(x^*)$ имеет обратную, откуда можно найти вектор-погрешность $\varepsilon^{(n+1)}$:

$$\varepsilon^{(n+1)} \approx - \frac{1}{2} [f'(x^*)]^{-1} \sum_{j,k=1}^v \frac{\partial^2 f(x^*)}{\partial x_j \partial x_k} \varepsilon_j^{(n)} \varepsilon_k^{(n)}. \quad (1.10.8)$$

Равенство говорит о том, что погрешности $\varepsilon_j^{(n+1)}$ приближения номера $n+1$ будут

малыми величинами второго порядка относительно погрешностей $\varepsilon_j^{(n)}$ предыдущего приближения номера n .

Это обстоятельство заставляет ожидать, что если в окрестности точного решения x^* определитель $D[f'(x)]$ матрицы Якоби отличен от нуля и если исходное приближение $x^{(0)}$ взято достаточно близким к x^* , ньютоновский процесс осуществим и последовательность $x^{(n)}$ будет сходиться к решению x^* по квадратичному закону.

Все сказанное ввиду приближенности равенства (1.10.8) может быть использовано лишь для ориентировочных, нестрогих оценок и заключений. Теоремы, выясняющие точные условия сходимости процесса и дающие строгие оценки погрешности, будут сформулированы ниже. Мы получим их как частные случаи из соответствующих теорем для операторных уравнений.

Система равенств

$$f_i(x_1, \dots, x_v) = y_i \quad (i = 1, 2, \dots, v).$$

или, коротко, $f(x) = y$ дает отображение v -мерного векторного пространства X на v -мерное векторное пространство Y . В каждом из пространств X и Y можно ввести свою норму. Для упрощения мы будем считать, что в X и Y введена одна и та же норма, и остановим свое внимание на нормах m , s и l (кубической, октаэдрической и шаровой).

Первой *) производной $f'(x)$ от оператора f на элементе x будет оператор линейного преобразования $X \rightarrow Y$

$$y_i' = \sum_{k=1}^v a_{ik} x_k', \quad (1.10.9)$$

матрицей которого является матрица Якоби с элементами

$$a_{ik} = \frac{\partial f_i(x)}{\partial x_k}.$$

Вторая производная $f''(x)$ на элементе x может быть истолкована как билинейный оператор $y'' = P(x)(x', x'')$, где

$$y_i'' = \sum_{j, k=1}^v \frac{\partial^2 f_i(x)}{\partial x_j \partial x_k} x_j' x_k'' = \sum_{j, k=1}^v a_{ijk}^{(i)} x_j' x_k''. \quad (1.10.10)$$

Предположим сначала, что в X и Y введена кубическая норма $\|x\|_m = \max_i |x_i|$.

Рассмотрим произвольную матрицу

$$A = \begin{bmatrix} a_{11} & \dots & a_{1v} \\ \cdot & \cdot & \cdot \\ a_{v1} & \dots & a_{vv} \end{bmatrix}.$$

Как известно,**) норма матрицы, подчиненная кубической норме вектора, есть

$$\|A\|_m = \max_i \sum_{j=1}^v |a_{ij}|. \quad (1.10.11)$$

*) Значения $f'(x)$ и $f''(x)$ объяснены и вычислены в добавлении I.

**) См. § 2.1. В гл. 2 и 3 приняты другие обозначения норм векторов и матриц, привычные в книгах по линейной алгебре, а именно, нормы $\|\dots\|_m$, $\|\dots\|_s$ и $\|\dots\|_l$ обозначают соответственно $\|\dots\|_1$, $\|\dots\|_{II}$ и $\|\dots\|_{III}$.

В теореме 1 § 1.7, которую мы хотим применить к системе (1.10.1), особую роль играет оператор $\Gamma_0 = [f'(x_0)]^{-1}$. В нашей задаче он является матрицей, обратной для $f'(x^{(0)})$. Ее m -норма имеет следующее значение. Предположим, что определитель матрицы Якоби $D(f'(x^{(0)})) = D$ отличен от нуля, и пусть D_{jk} есть алгебраическое дополнение элемента $\frac{\partial f_j(x_0)}{\partial x_k}$. Согласно (1.10.11),

$$\|\Gamma_0\|_m = \frac{1}{|D|} \max_k \sum_{j=1}^v |D_{jk}|. \quad (1.10.12)$$

Найдем оценку для нормы второй производной $f''(x)$:

$$\begin{aligned} |y_i''| &= \left| \sum_{j,k=1}^v a_{jk}^{(i)} x_j' x_k'' \right| \leq \sum_{j=1}^v \left| \sum_{k=1}^v a_{jk}^{(i)} x_j' \right| \max_k |x_k''| \leq \\ &\leq \sum_{j,k=1}^v |a_{jk}^{(i)}| \max_j |x_j'| \max_k |x_k''| = \sum_{j,k=1}^v |a_{jk}^{(i)}| \|x'\| \|x''\|. \end{aligned}$$

Отсюда мы получаем следующую оценку нормы $f''(x)$:

$$\|f''(x)\|_m \leq \max_i \sum_{j,k=1}^v |a_{jk}^{(i)}| = \max_i \sum_{j,k=1}^v \left| \frac{\partial^2 f_i(x)}{\partial x_j \partial x_k} \right|. \quad (1.10.13)$$

Соотношения (1.10.12) и (1.10.13) позволяют сформулировать в m -метрике теорему о сходимости процесса Ньютона для системы (1.10.1), являющуюся частным случаем теоремы 1 § 1.9, в следующем виде.

Теорема 1. Пусть выполнены условия:

1) функции $f_i(x) = f_i(x_1, \dots, x_v)$ определены и дважды непрерывно дифференцируемы в области

$$|x_i - x_i^{(0)}| \leq \delta \quad (i=1, 2, \dots, v), \quad (\alpha)$$

при этом для вторых производных в этой области выполнено неравенство

$$\sum_{j,k=1}^v \left| \frac{\partial^2 f_i(x)}{\partial x_j \partial x_k} \right| \leq K \quad (i=1, 2, \dots, v);$$

2) значения $x_1^{(0)}, \dots, x_v^{(0)}$ образуют приближенное решение системы (1.10.1) и для них выполняется неравенство

$$|f_i(x^{(0)})| \leq \eta;$$

3) матрица Якоби $f'(x)$ имеет в точке $x^{(0)}(x_1^{(0)}, \dots, x_v^{(0)})$ определитель $D = D(f'(x^{(0)}))$, отличный от нуля, и если D_{jk} есть алгебраическое дополнение элемента $\frac{\partial f_i}{\partial x_k}$, то верна оценка

$$\frac{1}{|D|} \sum_{j=1}^v |D_{jk}| \leq B \quad (k=1, 2, \dots, v);$$

4) для чисел K, η, B выполняется условие

$$h = B^2 K \eta \leq \frac{1}{2};$$

5) для δ верно неравенство

$$\frac{1 - \sqrt{1-2h}}{h} B \eta \leq \delta.$$

Тогда:

1) система (1.10.1) имеет в области (α) решение $x^* = (x_1^*, \dots, x_v^*)$;

2) последовательность ньютоновских приближений $x^{(n)} = (x_1^{(n)}, \dots, x_v^{(n)})$ может быть построена, принадлежит области (α) и сходится к решению x^* ;

3) скорость сходимости оценивается неравенством

$$|x_i^* - x_i^{(n)}| \leq t^* - t_n \quad (i=1, 2, \dots, v),$$

где $t^* = \frac{1 - \sqrt{1-2h}}{h} B \eta$ есть меньший корень квадратного уравнения $P(t) = \frac{1}{2} K t^2 - \frac{1}{B} t + \eta = 0$ и t_n ($n=0, 1, 2, \dots$) — последовательные приближения к нему по методу Ньютона, построенные при $t_0=0$.

Пусть теперь в пространствах X и Y введена октаэдрическая норма $\|x\|_* = \sum_{i=1}^v |x_i|$. Подчиненная ей норма матрицы (§ 2.1) есть

$$\|A\|_* = \max_j \sum_{i=1}^v |a_{ij}|. \quad (1.10.14)$$

Поэтому s -норма матрицы $\Gamma_0 = [f'(x^0)]^{-1}$ будет следующей:

$$\|\Gamma_0\|_* = \frac{1}{|D|} \max_j \sum_{k=1}^v |D_{jk}|, \quad (1.10.15)$$

где D_{jk} есть алгебраическое дополнение элемента $\frac{\partial f_j}{\partial x_k}$ в матрице $f'(x^{(0)})$.

Почти так же просто оценивается норма $y'' = f''(x)$.

$$\begin{aligned} |y_i''| &= \left| \sum_{j,k=1}^v a_{jk}^{(i)} x_j' x_k'' \right| \leq \sum_{j=1}^v |x_j'| \sum_{k=1}^v |a_{jk}^{(i)}| \cdot |x_k''| \leq \sum_{j=1}^v |x_j'| \max_k |a_{jk}^{(i)}| \times \\ &\times \sum_{k=1}^v |x_k''| \leq \max_{jk} |a_{jk}^{(i)}| \|x'\|_* \|x''\|_* \|y\|_* \leq \|x'\|_* \|x''\|_* \sum_i \max_{j,k} |a_{jk}^{(i)}|. \end{aligned}$$

Отсюда вытекает оценка s -нормы $f''(x)$.

$$\|f''(x)\|_s \leq \sum_{i=1}^v \max_{jk} |a_{jk}^{(i)}| = \sum_{i=1}^v \max_{jk} \left| \frac{\partial^2 f_i(x)}{\partial x_j \partial x_k} \right| \quad (1.10.16)$$

и мы можем сформулировать теорему в s -метрике об осуществимости и сходимости алгоритма Ньютона для системы (1.10.1), как частный случай теоремы 1 § 1.9.

Теорема 2. Пусть выполнены условия:

1) если функции $f_i(x) = f_i(x_1, \dots, x_v)$ определены и имеют непрерывные вторые производные в области

$$\sum_{i=1}^v |x_i - x_i^{(0)}| \leq \delta^1 \quad (\alpha)$$

и

$$\left| \frac{\partial^2 f_i(x)}{\partial x_j \partial x_k} \right| \leq L_i \quad (j, k = 1, 2, \dots, v; \quad x \in (\alpha); \quad i = 1, 2, \dots, v);$$

2) значения $x_1^{(0)}, \dots, x_v^{(0)}$ образуют приближенное решение системы (1.10.1) и выполняется неравенство

$$\sum_{i=1}^v |f_i(x^{(0)})| \leq \eta;$$

3) матрица Якоби $f'(x^{(0)})$ имеет определитель $D = D[f'(x^{(0)})]$, отличный от нуля, и верна оценка

$$\frac{1}{|D|} \sum_{k=1}^v |D_{jk}| \leq B \quad (j = 1, 2, \dots, v);$$

4) для чисел K, η, B выполняется условие

$$h = B^2 K \eta \leq \frac{1}{2}, \quad K = L_1 + L_2 + \dots + L_v;$$

5) для δ справедливо неравенство

$$\frac{1 - \sqrt{1 - 2h}}{h} B \eta \leq \delta.$$

Тогда:

- 1) система (1.10.1) в области (α) имеет решение x^* ;
- 2) ньютоновский процесс (1.10.5) для системы (1.10.1) осуществим, последовательность приближений $x^{(n)}$ лежит в области (α) и сходится к решению x^* ;
- 3) быстрота сходимости оценивается неравенством

$$\sum_{i=1}^v |x_i^* - x_i^{(n)}| \leq t^* - t_n,$$

где t^* и t_n определены, как в теореме 1.

Немного более сложным является случай (§ 2.1), когда в пространства X и Y введена шаровая метрика:

$$\|x\|_l = \left[\sum_{i=1}^v x_i^2 \right]^{\frac{1}{2}}.$$

Подчиненная ей норма действительной матрицы $A = [a_{jk}]$ равна корню квадратному из наибольшего собственного значения Λ произведения $A'A$, где A' есть транспонированная матрица:

$$\|A\|_l = \sqrt{\Lambda}. \quad (1.10.17)$$

Находить Λ часто затруднительно, и в приложениях предпочитают иногда пользоваться не точным значением, а оценкой нормы матрицы, которая легко может быть получена.

$$|y_j''|^2 = \left| \sum_{k=1}^v a_{jk} x_k \right|^2 \leq \sum_{k=1}^v a_{jk}^2 \sum_{k=1}^v x_k^2 = \sum_{k=1}^v a_{jk}^2 \|x\|_l^2,$$

$$\|y''\|_l^2 \leq \sum_{j,k=1}^v a_{jk}^2 \|x\|_l^2,$$

$$\|A\|_l \leq \left\{ \sum_{j,k=1}^v a_{jk}^2 \right\}^{\frac{1}{2}}. \quad (1.10.18)$$

Рассмотрим теперь матрицу $\Gamma_0 = [f'(x^{(0)})]^{-1}$. Оценка (1.10.18) ее нормы будет, очевидно, следующей:

$$\|\Gamma_0\|_l \leq \frac{1}{|D|} \left\{ \sum_{j,k=1}^v D_{jk}^2 \right\}^{\frac{1}{2}}. \quad (1.10.19)$$

Получим, наконец, оценку для нормы $f''(x)$. Неравенство Коши — Буняковского дает

$$(y_i'')^2 = \left| \sum_{j,k=1}^v a_{jk}^{(i)} x_j' x_k'' \right|^2 \leq \left\{ \sum_{j=1}^v x_j'^2 \right\} \sum_{k=1}^v \left[\sum_{j=1}^v a_{jk}^{(i)} x_k'' \right]^2 \leq \Lambda_i \|x'\|_l^2 \|x''\|_l^2,$$

где Λ_i означает наибольшее собственное значение матрицы $A_i' A_i$, являющейся произведением матрицы

$$A_i = \begin{bmatrix} a_{11}^{(i)} & \cdots & a_{1v}^{(i)} \\ \vdots & \ddots & \vdots \\ a_{v1}^{(i)} & \cdots & a_{vv}^{(i)} \end{bmatrix}, \quad a_{jk}^{(i)} = \frac{\partial^2 f_i(x)}{\partial x_j \partial x_k},$$

и транспонированной A_i' . Отсюда получаем

$$\|y''\|_l^2 = \sum_{i=1}^v (y_i'')^2 \leq \sum_{i=1}^v \Lambda_i \|x'\|_l^2 \|x''\|_l^2.$$

Это дает оценку нормы $f''(x)$:

$$\|f''(x)\|_l \leq \left\{ \sum_{i=1}^v \Lambda_i \right\}^{\frac{1}{2}}. \quad (1.10.20)$$

В нашей задаче пользоваться этой оценкой затруднительно, так как пришлось бы находить наибольшие собственные значения Λ_i для всяких $x = (x_1^*, \dots, x_v)$ из области (а). Предпочтительнее воспользоваться другой, более простой, но несколько более грубой оценкой. Для ее получения достаточно обратить внимание на то, что из сравнения соотношений (1.10.17) и (1.10.18) сразу следует $\Lambda_i \leq \sum_{j,k} a_{jk}^{(i)2}$ и

$$\|f''(x)\| \leq \left\{ \sum_{i,j,k=1}^v a_{jk}^{(i)2} \right\}^{\frac{1}{2}} = \left\{ \sum_{i,j,k=1}^v \left[\frac{\partial^2 f_i(x)}{\partial x_j \partial x_k} \right]^2 \right\}^{\frac{1}{2}}. \quad (1.10.21)$$

Неравенства (1.10.19) и (1.10.20) позволяют сформулировать в l -норме теорему о сходимости процесса Ньютона (1.10.5) для системы (1.10.1), как частный случай теоремы 1 § 1.9.

Теорема 3. Пусть выполняются условия:

1) функции $f_i(x) = f_i(x_1, \dots, x_v)$ определены и дважды непрерывно дифференцируемы в области

$$\sum_{i=1}^v (x_i - x_i^{(0)})^2 \leq \delta^2, \quad (\alpha)$$

при этом для вторых производных в этой области выполняется неравенство

$$\sum_{i,j,k=1}^v \left[\frac{\partial^2 f_i(x)}{\partial x_j \partial x_k} \right]^2 \leq K^2;$$

2) для исходного приближения $x^{(0)}(x_1^{(0)}, \dots, x_v^{(0)})$ к решению системы справедливо неравенство

$$\sum_{i=1}^v f_i^2(x^{(0)}) \leq \eta^2;$$

3) матрица Якоби $f'(x^{(0)})$ имеет определитель $D = D[f'(x^{(0)})]$, отличный от нуля, и верна оценка

$$\frac{1}{|D|} \left\{ \sum_{j,k=1}^v D_{jk}^2 \right\}^{\frac{1}{2}} \leq B.$$

4) для чисел K, η, B выполняется условие

$$h = B^2 K \eta \leq \frac{1}{2};$$

5) для δ верно неравенство

$$\frac{1 - \sqrt{1 - 2h}}{h} B \eta \leq \delta.$$

Тогда:

- 1) система (1.10.1) имеет в области (α) решение;
- 2) ньютоновский процесс (1.10.5) для системы (1.10.1) осуществим, последовательность приближений $x^{(n)}$ лежит в области (α) и сходится к решению x^* ;
- 3) скорость сходимости может быть оценена неравенством

$$\left\{ \sum_{i=1}^v (x_i^* - x_i^{(n)})^2 \right\}^{\frac{1}{2}} \leq t^* - t_n.$$

Значения t^* и t_n указаны в теореме 1.

§ 1.11. МЕТОД РЕШЕНИЯ, ОСНОВАННЫЙ НА ВОЗВЕДЕНИИ КОРНЕЙ В СТЕПЕНЬ

Этот метод был предложен независимо друг от друга несколькими учеными. В нашей научной и учебной литературе он заслуженно называется методом Лобачевского. Применяют его для нахождения решений алгебраических уравнений, хотя, если иметь в виду принципиальную сторону дела, он может быть применен для уравнений в аналитических функциях. Как будет видно из дальнейшего, метод Лобачевского не требует предварительного приближенного нахождения корней и позволяет одновременно найти все корни многочлена. Недостатком его является тот факт, что при вычислениях приходится иметь дело с числами, сильно различающимися по порядкам величин.

Будем рассматривать алгебраическое уравнение степени n

$$P(x) = a_0 x^n + a_1 x^{n-1} + \dots + a_n = 0. \quad (1.11.1)$$

Корни его перенумеруем в порядке убывания модулей:

$$|x_1| \geq |x_2| \geq |x_3| \geq \dots \geq |x_n|.$$

В основание метода был положен весьма простой факт. Напомним соотношения между корнями и коэффициентами многочлена:

$$\left. \begin{aligned} x_1 + x_2 + \dots + x_n &= -\frac{a_1}{a_0}, \\ x_1x_2 + x_1x_3 + \dots + x_{n-1}x_n &= \frac{a_2}{a_0}, \\ x_1x_2x_3 + \dots + x_{n-2}x_{n-1}x_n &= -\frac{a_3}{a_0}, \\ &\vdots \\ x_1x_2 \dots x_n &= (-1)^n \frac{a_n}{a_0}. \end{aligned} \right\} \quad (1.11.2)$$

Условимся говорить, что корни x_i ($i=1, 2, \dots, n$) сильно разделены в смысле отношений, если модуль предыдущего корня будет во много раз больше модуля следующего корня. Записать это можно в форме неравенства

$$\left| \frac{x_k}{x_{k+1}} \right| \gg 1 \quad (k=1, 2, \dots, n-1).$$

Так как в настоящем параграфе мы будем говорить о сильной разделенности лишь в смысле отношений, то для сокращения фраз слова «в смысле отношений» будем опускать и говорить только «корни сильно разделены».

Если имеет место сильная разделенность корней, то равенства Виета (1.11.2) упрощаются. В самом деле, если в первом из них вынести за скобки x_1 , оно примет форму

$$x_1 \left(1 + \frac{x_2}{x_1} + \dots + \frac{x_n}{x_1} \right) = -\frac{a_1}{a_0}.$$

Когда разделенность корней является достаточно сильной, то все отношения, стоящие в скобках, будут величинами, пренебрежимо малыми сравнительно с единицей. В пределах принятой точности их можно будет отбросить и заменить первое равенство следующим: $x_1 \approx -\frac{a_1}{a_0}$. Аналогич-

ное будет иметь место и для всех других равенств Виета, и (1.11.2) заменится следующей системой приближенных равенств, верных лишь в принятой точности вычислений:

$$\left. \begin{aligned} x_1 &\approx -\frac{a_1}{a_0}, \\ x_1x_2 &\approx \frac{a_2}{a_0}, \\ x_1x_2x_3 &\approx -\frac{a_3}{a_0}, \\ . & \\ x_1x_2 \dots x_n &\approx (-1)^n \frac{a_n}{a_0}. \end{aligned} \right\} \quad (1.11.3)$$

До сих пор словам «корни сильно разделены» мы не придавали строгого арифметического смысла и пользовались ими только как средством наглядного описания. Сейчас можно сказать, что в нашем изложении корни называются сильно разделенными, если в равенствах Виета первые члены в левых частях будут главными и сумма модулей остальных членов в каждой из левых частей будет лежать вне принятой точности вычислений и ею можно пренебречь.

Если корни сильно разделены, то они очень просто находятся по коэффициентам. Из (1.11.3) следует, что

$$x_1 \approx -\frac{a_1}{a_0}, \quad x_2 \approx -\frac{a_2}{a_1}, \quad x_3 \approx -\frac{a_3}{a_2}, \quad \dots, \quad x_n \approx -\frac{a_n}{a_{n-1}}.$$

Изложенные простые соображения подсказывают путь для разыскания корней. Если в заданном уравнении корни не разделены сильно, то мы можем надеяться добиться их сильного разделения, если возведем их в высокие степени. Чтобы осуществить это, нужно указать достаточно простой алгоритм, позволяющий по заданному уравнению построить новое уравнение, корни которого были бы степенями его корней. Достаточно, очевидно, построить алгоритм нахождения уравнения, корни которого будут квадратами корней заданного уравнения. Выполняя несколько шагов преобразований, мы будем возводить корни в степени 2, 4, 8, ..., $m=2^k, \dots$

Наряду с

$$P(x) = a_0(x-x_1)(x-x_2) \dots (x-x_n)$$

рассмотрим многочлен $P^*(x)$ с тем же коэффициентом a_0 при старшей степени, корни которого равны $-x_1, -x_2, \dots, -x_n$,

$$P^*(x) = a_0(x+x_1)(x+x_2) \dots (x+x_n).$$

Чтобы получить многочлен $P_1(y)$ с корнями $x_1^2, x_2^2, \dots, x_n^2$, достаточно в произведении PP^* заменить x^2 на y :

$$P(x)P^*(x) = a_0^2(x^2-x_1^2) \dots (x^2-x_n^2) = a_0^2(y-x_1^2) \dots (y-x_n^2) = P_1(y).$$

Многочлен $P^*(x)$ строится просто. Равенства (1.11.2) показывают, что при сохранении a_0 , когда корни x_1, \dots, x_n заменяются на $-x_1, \dots, -x_n$, коэффициенты с нечетными номерами a_1, a_3, \dots переходят в $-a_1, -a_3, \dots$, а коэффициенты с четными номерами сохраняют значения:

$$P^*(x) = a_0x^n - a_1x^{n-1} + a_2x^{n-2} - \dots + (-1)^n a_n.$$

Запишем $P_1(y)$ в форме

$$P_1(y) = a_0^{(1)}y^n + a_1^{(1)}y^{n-1} + \dots + a_n^{(1)}.$$

Перемножение многочленов

$$P(x)P^*(x) = (a_0x^n + a_1x^{n-1} + a_2x^{n-2} + \dots)(a_0x^n - a_1x^{n-1} - a_2x^{n-2} - \dots)$$

приводит к следующему правилу вычисления коэффициентов $a_k^{(1)}$:

$$\left. \begin{aligned} a_0^{(1)} &= a_0^2, \\ a_1^{(1)} &= 2a_0a_2 - a_1^2, \\ a_2^{(1)} &= 2a_0a_4 - 2a_1a_3 + a_2^2, \\ a_3^{(1)} &= 2a_0a_6 - 2a_1a_5 + 2a_2a_4 - a_3^2, \\ &\dots \\ a_n^{(1)} &= (-1)^n a_n^2. \end{aligned} \right\} \quad (1.11.4)$$

Применяя это преобразование к многочлену $P_1(y)$, построим многочлен $P_2(y)$ с корнями x_1^4, \dots, x_n^4 и т. д. После k -кратного преобразования получим многочлен

$$P_k(x) = a_0^{(k)}x^n + a_1^{(k)}x^{n-1} + \dots + a_n^{(k)},$$

корнями которого будут x_1^m, \dots, x_n^m при $m=2^k$. Соотношения вида (1.11.2) между корнями и коэффициентами для него будут

$$\left. \begin{aligned} x_1^m + x_2^m + \dots + x_n^m &= -\frac{a_1^{(k)}}{a_0^{(k)}}, \\ x_1^m x_2^m + \dots + x_{n-1}^m x_n^m &= \frac{a_2^{(k)}}{a_0^{(k)}}, \\ x_1^m x_2^m x_3^m + \dots + x_{n-2}^m x_{n-1}^m x_n^m &= -\frac{a_3^{(k)}}{a_0^{(k)}}, \\ &\dots \\ x_1^m x_2^m \dots x_n^m &= (-1)^n \frac{a_n^{(k)}}{a_0^{(k)}}. \end{aligned} \right\} \quad (1.11.5)$$

Нашей задачей сейчас будет установить правила, следуя которым можно будет получать при больших m из равенств (1.11.5) сведения о численных значениях корней x_i .

Исследование общего случая уравнения связано с громоздкой записью и сложными рассуждениями, которые могут затруднить пони-

мание сравнительно простого содержания задачи. Мы ограничимся поэтому разбором простых типичных случаев, понимание которых будет достаточным для многих приложений. Кроме того, они позволяют составить представление об использовании (1.11.5) в более общих случаях.

В последующем изложении будем считать коэффициенты a_k уравнения (1.11.1) действительными.

1. Предположим, что корни уравнения действительны и различны по абсолютной величине:

$$|x_1| > |x_2| > \dots > |x_n|. \quad (1.11.6)$$

Сделав k достаточно большим, можно добиться того, чтобы первые члены в левых частях равенств (1.11.5) были бы преобладающими и суммы остальных членов лежали бы вне принятой точности. Тогда вместо (1.11.5) мы получим приближенные равенства, верные на принятое число значащих цифр:

$$x_1^m \approx -\frac{a_1^{(k)}}{a_0^{(k)}}, \quad x_1^m x_2^m \approx -\frac{a_2^{(k)}}{a_0^{(k)}}, \quad x_1^m x_2^m x_3^m \approx -\frac{a_3^{(k)}}{a_0^{(k)}}, \dots \quad (1.11.7)$$

Из них находим

$$x_1^m \approx -\frac{a_1^{(k)}}{a_0^{(k)}}, \quad x_2^m \approx -\frac{a_2^{(k)}}{a_1^{(k)}}, \quad x_3^m \approx -\frac{a_3^{(k)}}{a_2^{(k)}}, \dots \quad (1.11.8)$$

Отсюда мы получим абсолютные значения корней $|x_1|$, $|x_2|$, $|x_3|$, ... Знаки же их можно определить, например, при помощи подстановки в уравнение.

При преобразованиях мы вычисляем коэффициенты $a_i^{(k)}$ многочленов $P_k(x)$. Нам осталось еще указать, как можно по поведению $a_i^{(k)}$ судить о том, достигли мы или не достигли необходимого k . Допустим, что нужное значение k уже достигнуто и с принятой точностью равенства (1.11.7) выполняются. Прделаем еще одно, по сути дела уже лишнее, преобразование и найдем многочлен

$$P_{k+1}(x) = a_0^{(k+1)}x^n + a_1^{(k+1)}x^{n-1} + \dots$$

Для него также должны быть верными равенства вида (1.11.7)

$$x_1^{2m} \approx -\frac{a_1^{(k+1)}}{a_0^{(k+1)}}, \quad x_1^{2m} x_2^{2m} \approx -\frac{a_2^{(k+1)}}{a_0^{(k+1)}}, \dots$$

Так как $a_0^{(k+1)} = [a_0^{(k)}]^2$, из сравнения этих равенств с (1.11.7) следует, что

абсолютные величины коэффициентов $a_i^{(k+1)}$ должны быть в принятой точности равны квадратам коэффициентов $a_i^{(k)}$:

$$|a_1^{(k+1)}| = [a_1^{(k)}]^2, \quad |a_2^{(k+1)}| = [a_2^{(k)}]^2, \dots$$

Выполнение этих равенств и будет свидетельствовать о том, что необходимое значение k уже было достигнуто на предпоследнем шаге и вычисление коэффициентов $a_i^{(k+1)}$ многочлена $P_{k+1}(x)$ было лишним.

Для расчетной схемы этот признак означает следующее: коэффициенты $a_i^{(h)}$ и $a_i^{(h+1)}$ связаны между собой равенствами вида (1.11.4), нужно в них только a_i заменить на $a_i^{(h)}$ и $a_i^{(4)}$ — на $a_i^{(h+1)}$. Вычисления следует прекратить, если в принятой точности в правых частях сохранятся только квадраты коэффициентов, а суммы удвоенных произведений окажутся ниже границы точности.

2. Допустим теперь, что корни действительные и среди них есть равные по абсолютной величине. Простоты ради предположим, что только два корня, например x_2 и x_3 , имеют одинаковое абсолютное значение:

$$|x_1| > |x_2| = |x_3| > |x_4| > \dots > |x_n|.$$

При достаточно больших k соотношения (1.11.5) примут следующую приближенную форму:

$$\left. \begin{aligned} x_1^m &\approx -\frac{a_1^{(k)}}{a_0^{(k)}}, \\ 2x_1^mx_2^m &\approx \frac{a_2^{(k)}}{a_0^{(k)}}, \\ x_1^mx_2^{2m} &\approx -\frac{a_3^{(k)}}{a_0^{(k)}}, \\ x_1^mx_2^mx_4^m &\approx \frac{a_4^{(k)}}{a_0^{(k)}}, \\ . & \\ x_1^mx_2^mx_4^m\dots x_n^m &= (-1)^n \frac{a_n^{(k)}}{a_0^{(k)}}. \end{aligned} \right\} \quad (1.11.9)$$

Для нахождения x_2^m можно воспользоваться с одинаковым правом вто-

рым и третьим равенством. Например, из первого и третьего равенства следует $x_2^{2m} \approx \frac{a_3^{(k)}}{a_1^{(k)}}$. Далее последовательно находим $x_4^m \approx -\frac{a_4^{(k)}}{a_3^{(k)}}, \dots$

Нужно вновь указать признак, как можно по поведению коэффициентов $a_i^{(k)}$ судить о том, что необходимое значение k уже достигнуто. Допустим, что с принятой точностью равенства (1.11.9) верны. Сделаем еще раз преобразование и от многочлена $P_k(x)$ перейдем к

$$P_{k+1}(x) = a_0^{(k+1)}x^n + a_1^{(k+1)}x^{n-1} + \dots$$

Соотношения (1.11.9) для него будут

$$x_1^{2m} \approx -\frac{a_1^{(k+1)}}{a_0^{(k+1)}}, \quad 2x_1^{2m}x_2^{2m} \approx \frac{a_2^{(k+1)}}{a_0^{(k+1)}}, \quad x_1^{2m}x_2^{4m} \approx -\frac{a_3^{(k+1)}}{a_0^{(k+1)}}, \dots$$

Сравнение их с (1.11.9), ввиду $a_0^{(k+1)} = [a_0^{(k)}]^2$, приводит к равенствам

$$a_0^{(k+1)} = [a_0^{(k)}]^2, \quad |a_1^{(k+1)}| = |a_1^{(k)}|^2, \quad a_2^{(k+1)} = \frac{1}{2} [a_2^{(k)}]^2, \quad |a_3^{(k+1)}| = [a_3^{(k)}]^2, \dots$$

Как видно, все коэффициенты $a_i^{(k+1)}$ по абсолютной величине равны квадратам коэффициентов $a_i^{(k)}$, кроме $a_2^{(k+1)}$, для которого будет:

$$a_2^{(k+1)} = \frac{1}{2} [a_2^{(k)}]^2.$$

Выполнение таких равенств и является признаком того, что достаточно большое для нашей цели значение k достигнуто на предпоследнем шаге и переход от k к $k+1$ был излишним.

3. Пусть уравнение имеет комплексные корни. Они попарно сопряжены, так как a_i ($i=0, 1, \dots, n$) действительны. Предположим, что существует только одна пара комплексных корней, и пусть это будут

$$x_2 = re^{i\varphi} \quad \text{и} \quad x_3 = re^{-i\varphi}.$$

$$|x_1| > r > |x_4| > \dots > |x_n|.$$

Ввиду $x_2^m + x_3^m = 2r^m \cos m\varphi$ и $x_2^m x_3^m = r^{2m}$, соотношения (1.11.5) будут такими:

• • • • •

лениях число правильных знаков:

• • • • • • • • • •

Из них находим

$$x_1^m \approx -\frac{a_1}{a_0^{(k)}}, \quad r^{2m} \approx -\frac{a_3}{a_1^{(k)}}, \quad x_4^m \approx -\frac{a_4}{a_3^{(k)}}, \dots$$

и после этого можем вычислить абсолютные значения действительных корней и модуль r двух комплексных корней x_2 и x_3 . Знаки корней x_1, x_4, \dots, x_n могут быть определены при помощи подстановки в уравнение.

Для нахождения же x_2 и x_3 можно воспользоваться, например, равенством $x_1 + x_2 + x_3 + \dots + x_n = -\frac{a_1}{a_0}$. В нашем случае оно даст

$$x_1 + 2r \cos \varphi + x_4 + \dots + x_n = -\frac{a_1}{a_0}$$

и позволит вычислить $\cos \varphi$. Затем найдем $\sin \varphi = (1 - \cos^2 \varphi)^{1/2}$ и корни $x_2 = \bar{x}_3 = r(\cos \varphi + i \sin \varphi)$.

По причинам, которые выяснились в предыдущих двух случаях, преобразования могут быть остановлены на таком шаге k , чтобы при переходе от многочлена $P_k(x)$ к $P_{k+1}(x)$ были на принятое число знаков выполнены равенства

$$a_0^{(k+1)} = [a_0^{(k)}]^2, \quad |a_1^{(k+1)}| = [a_1^{(k)}]^2, \dots, \quad |a_n^{(k+1)}| = [a_n^{(k)}]^2.$$

При росте k в поведении коэффициента $a_2^{(k)}$, ввиду наличия в первом члене левой части второго равенства члена $2x_1^m r^m \cos \varphi$, не будет регулярности. Коэффициент $a_2^{(k)}$ будет, вообще говоря, колебаться по абсолютной величине и изменять свой знак. Этот факт будет указывать на то, что у многочлена $P(x)$ корни x_2 и x_3 являются комплексными.

§ 1.12. НАХОЖДЕНИЕ КОРНЕЙ МНОГОЧЛЕНОВ ПРИ ПОМОЩИ ВЫДЕЛЕНИЯ МНОЖИТЕЛЕЙ

Пусть дано алгебраическое уравнение

$$P(x) = a_0 x^n + a_1 x^{n-1} + \dots + a_n = 0 \quad (a_n \neq 0). \quad (1.12.1)$$

Чем выше степень n уравнения, тем, вообще говоря, труднее задача его решения. Разложение многочлена $P(x)$ на множители позволяет часто весьма сильно упростить задачу нахождения корней. Существуют алгоритмы, позволяющие выделить из $P(x)$ множитель любой наперед заданной степени m ($1 \leq m < n$).

Мы остановимся лишь на простейших задачах такого разложения. Будем предполагать, что коэффициенты a_k ($k=0, 1, \dots, n$) есть действительные числа. Алгоритмы, о которых мы будем говорить ниже, при некоторых изменениях легко могут быть перенесены на многочлены с комплексными коэффициентами.

Известно, что всякий действительный многочлен $P(x)$ может быть разложен на действительные множители: линейные, отвечающие действительным корням многочлена $P(x)$, и квадратные, отвечающие парам сопряженных комплексных корней. Это позволяет ограничиться изучением задачи выделения из $P(x)$ множителей только первой и второй сте-

пеней. Принципиально говоря, можно было бы ограничиться, как это часто делают, выделением лишь квадратных множителей. Но мы не будем исключать задачу выделения линейных множителей, так как вычисления для их нахождения несколько проще, чем для нахождения множителей второй степени.

Начнем с задачи выделения линейного множителя и рассмотрим алгоритм предпоследнего остатка, предложенный Лином.

Пусть нахождению подлежит действительный корень α уравнения (1.12.1) и нужно выделить множитель $x - \alpha$, отвечающий этому корню. Предположим, что мы знаем приближение x_0 к α и можем составить лишь приближенное значение $x - x_0$ множителя $x - \alpha$. Для улучшения его будем делить $P(x)$ по обычным алгебраическим правилам на $x - x_0$. Если деление выполнить до конца, то в остатке получится постоянная величина. Остановимся на предпоследнем шаге. Предпоследний остаток будет линейным, вида $d_0x + d_1$. Предполагая $d_0 \neq 0$, запишем его в форме $d_0(x - x_1)$. Разность $x - x_1$ часто называют приведенным предпоследним остатком. Положим

$$P(x) = (x - x_0)(b_0^{(1)}x^{n-1} + b_1^{(1)}x^{n-2} + \dots + b_{n-2}^{(1)}x) + d_0(x - x_1),$$

x_1 примем за первое «улучшенное» значение корня.

Из последнего тождества найдем

$$P(x_0) = d_0(x_0 - x_1), \quad P(0) = -d_0x_1.$$

Исключая отсюда d_0 , получим явное выражение x_1 через исходное приближение:

$$x_1 = \frac{P(0)x_0}{P(0) - P(x_0)} = - \frac{a_n}{a_0x_0^{n-1} + a_1x_0^{n-2} + \dots + a_{n-1}}.$$

Для получения x_2 делим $P(x)$ на $x - x_1$, вновь находим «предпоследний остаток» и представляем его в приведенной форме $d_1(x - x_2)$ и т. д.

Если известно приближение x_k , то следующее приближение x_{k+1} находится по тому же алгоритму путем деления $P(x)$ на $x - x_k$ до предпоследнего остатка и приведения последнего к виду $d_k(x - x_{k+1})$:

$$P(x) = (x - x_k)(c_0x^{n-1} + c_1x^{n-2} + \dots + c_{n-2}x) + d_k(x - x_{k+1}).$$

Условием возможности алгоритма будет неравенство $d_k \neq 0$ ($k=0, 1, 2, \dots$).

Явное выражение x_{k+1} через x_k имеет вид

$$x_{k+1} = \frac{P(0)x_k}{P(0) - P(x_k)} = - \frac{a_n}{a_0x_k^{n-1} + a_1x_k^{n-2} + \dots + a_{n-1}}.$$

По существу имеем здесь дело с простым одношаговым итерационным процессом для нахождения решения уравнения

$$\begin{aligned} x = \varphi(x) &= x + \frac{xP(x)}{P(0) - P(x)} = \frac{xP(0)}{P(0) - P(x)} = \\ &= - \frac{a_n}{a_0x^{n-1} + a_1x^{n-2} + \dots + a_{n-1}} \end{aligned}$$

с определенным алгоритмом вычисления значения правой части уравнения в точке $x = x_h$. Сходимость процесса зависит от значений $\varphi'(x)$ вблизи решения $x = \alpha$. Но

$$\varphi'(\alpha) = 1 + \alpha \frac{P'(\alpha)}{P(\alpha)}.$$

Если

$$|\varphi'(\alpha)| = \left| 1 + \alpha \frac{P'(\alpha)}{P(\alpha)} \right| < 1,$$

то найдется такая окрестность решения $|x - \alpha| \leq q < 1$, в которой $|\varphi'(x)| \leq q < 1$. Если x_0 взято из этой окрестности, то можно ожидать, что последовательность x_n будет сходиться к решению α : $\lim x_n = \alpha$.

Отметим также, что из итерационного правила, которому можно придать форму

$$x_{h+1} = x_h + \frac{x_h}{P(0) - P(x_h)} P(x_h),$$

следует, что если x_h стремится к конечному пределу α , то для α должно выполняться равенство $P(\alpha) = 0$, т. е. α должен быть корнем уравнения.

Перейдем теперь к задаче выделения множителя второй степени. Мы будем искать его в форме $x^2 + px + q$ и предположим, что для коэффициентов его p и q указаны каким-либо путем приближенные значения p_0 и q_0 .

Выполним деление $P(x)$ на $x^2 + p_0x + q_0$. Если эту операцию выполнить до конца, то в остатке получится, вообще говоря, многочлен первой степени. Мы остановимся на предпоследнем остатке, который, как правило, будет многочленом второй степени $ax^2 + bx + c$. Разделим его на a , предполагая $a \neq 0$, и преобразуем остаток к приведенной форме $x^2 + p_1x + q_1$. Путем такой же операции деления $P(x)$ на $x^2 + p_1x + q_1$ строим второй предпоследний приведенный остаток и т. д. Мы получим, вообще говоря, последовательность приведенных предпоследних остатков $x^2 + p_kx + q_k$. Если окажется, что p_k и q_k имеют конечные пределы:

$$\lim_{k \rightarrow \infty} p_k = p^* \quad \text{и} \quad \lim_{k \rightarrow \infty} q_k = q^*,$$

то, как будет видно из дальнейшего изложения, многочлен x^2+px+q будет делителем $P(x)$.

Рассмотрим более внимательно алгоритм деления многочлена $P(x)$ на трехчлен x^2+px+q . Если выполнить деление до конца, получим равенство

$$\begin{aligned} P(x) &= a_0x^n + a_1x^{n-1} + \dots = \\ &= (x^2+px+q)(b_0x^{n-2} + b_1x^{n-3} + \dots + b_{n-3}x + b_{n-2}) + b_{n-1}(x+p) + b_n. \end{aligned}$$

Остаток от деления записан в особой форме, позволяющей просто и единообразно записать уравнения для определения b_i ($i=0, 1, \dots, n$).

Сравнение коэффициентов при степенях x даст равенства

$$\left. \begin{aligned} a_0 &= b_0, \\ a_1 &= pb_0 + b_1, \\ a_2 &= qb_0 + pb_1 + b_2, \\ &\dots \dots \dots \\ a_{n-1} &= qb_{n-3} + pb_{n-2} + b_{n-1}, \\ a_n &= qb_{n-2} + pb_{n-1} + b_n. \end{aligned} \right\} \quad (1.12.2)$$

Из них последовательно могут быть найдены b_0, b_1, \dots , при этом ясно, что b_i будут многочленами от p и q . Степени их относительно p и q легко подсчитать и установить, что b_i есть многочлен степени i относительно p и степени $i-1$ относительно q . Ниже, если нужно будет указать на зависимость b_i от p и q , мы будем их обозначать $b_i(p, q)$.

Попутно заметим, что x^2+px+q будет делителем $P(x)$ в том случае, когда остаток $b_{n-1}(x+p) + b_n$ будет тождественным нулем, что равносильно выполнению системы

$$b_n(p, q) + pb_{n-1}(p, q) = 0, \quad b_{n-1}(p, q) = 0 \quad (1.12.3)$$

или, что то же самое,

$$b_n(p, q) = 0, \quad b_{n-1}(p, q) = 0. \quad (1.12.4)$$

Предпоследний остаток от деления $P(x)$ на x^2+px+q есть

$$b_{n-2}(x^2+px+q) + b_{n-1}(x+p) + b_n = b_{n-2}x^2 + (pb_{n-2} + b_{n-1})x + a_n.$$

В описанном выше алгоритме нахождения приведенных предпоследних остатков следующие значения коэффициентов p_{k+1}, q_{k+1} находятся по предыдущим p_k, q_k по правилу

$$p_{k+1} = p_k + \frac{b_{n-1}(p_k, q_k)}{b_{n-2}(p_k, q_k)}, \quad q_{k+1} = \frac{a_n}{b_{n-2}(p_k, q_k)}. \quad (1.12.5)$$

Условием осуществимости алгоритма является неравенство

$$b_{n-2}(p_k, q_k) \neq 0 \quad (k=0, 1, 2, \dots).$$

Правила (1.12.5) весьма просто связаны с системой (1.12.3). Если заметить, что $a_n = qb_{n-2} + pb_{n-1} + b_n$ и, стало быть, второе равенство (1.12.5) равносильно

$$q_{k+1} = q_k + \frac{b_n(p_k, q_k) + p_k b_{n-1}(p_k, q_k)}{b_{n-2}(p_k, q_k)},$$

и записать систему в виде

$$\begin{aligned} p &= \varphi(p, q) = p + \frac{b_{n-1}(p, q)}{b_{n-2}(p, q)}, \\ q &= \psi(p, q) = q + \frac{b_n(p, q) + p b_{n-1}(p, q)}{b_{n-2}(p, q)}, \end{aligned} \quad (1.12.6)$$

становится ясным, что правило (1.12.5) есть не что иное, как простой одношаговый итерационный процесс для системы (1.12.3), преобразованной к форме (1.12.6).

Закончим изложение метода Лина доказательством факта, на который мы обращали внимание выше: если p_k и q_k имеют конечные пределы $\lim p_k = p^*$ и $\lim q_k = q^*$, то трехчлен $x^2 + p^*x + q^*$ будет делителем $P(x)$. Для этого достаточно обратиться к соотношению, связывающему два приближения $x^2 + p_k x + q_k$ и $x^2 + p_{k+1} x + q_{k+1}$:

$$\begin{aligned} P(x) &= (x^2 + p_k x + q_k) [b_0 x^{n-2} + b_1(p_k, q_k) x^{n-3} + \dots + b_{n-3}(p_k, q_k) x] + \\ &\quad + b_{n-2}(p_k, q_k) (x^2 + p_{k+1} x + q_{k+1}), \end{aligned}$$

где $b_i(p, q)$ есть многочлены от p и q . Если перейти здесь к пределу при $k \rightarrow \infty$, получим равенство, доказывающее утверждение:

$$\begin{aligned} P(x) &= (x^2 + p^* x + q^*) [b_0(p^*, q^*) x^{n-2} + b_1(p^*, q^*) x^{n-3} + \dots + \\ &\quad + b_{n-3}(p^*, q^*) x] + b_{n-2}(p^*, q^*) (x^2 + p^* x + q^*). \end{aligned}$$

Метод предпоследнего остатка сходится не во всех случаях и может расходиться, даже если исходные приближения p_0, q_0 взяты близко к точным значениям p и q . Кроме того, его сходимость может быть медленной. Часто указанных недостатков можно избежать, если для решения системы (1.12.4), к которой приводится нахождение квадратного делителя $P(x)$, применить другие методы решения. В частности, к решению системы можно применить метод Ньютона, или одно из его видоизменений. В этих последних методах в вычислениях для системы (1.12.4) придется

пользоваться значениями $b_n(p, q)$ и $b_{n-1}(p, q)$ и частных производных от них по p, q . Полезно обратить внимание на то, что значения b_n и b_{n-1} могут быть найдены без знания их явных выражений через p и q либо путем деления $P(x)$ на x^2+px+q с численными значениями p и q , либо при помощи рекурсионных уравнений (1.12.2). Оказывается, что и частные производные также могут быть найдены либо при помощи деления, либо из уравнений, сходных с (1.12.2), без знания явной зависимости b_n и b_{n-1} от p и q .

Напомним, что если $P(x)$ делить на x^2+px+q , то мы получим следующее его выражение

$$P(x) = (x^2+px+q)Q(x) + b_{n-1}(x+p) + b_n,$$

$$Q(x) = b_0x^{n-2} + b_1x^{n-3} + \dots + b_{n-2},$$

при этом b_i вычисляются либо путем деления, либо при помощи системы равенств (1.12.2), которую коротко можно записать так:

$$b_j = a_j - pb_{j-1} - qb_{j-2} \quad (j=1, 2, \dots, n), \quad (1.12.7)$$

$$b_{-1} = b_{-2} = 0.$$

Для нахождения частных производных от b_j вычислим производные по p и q от обеих частей уравнения (1.12.7):

$$\frac{\partial}{\partial p} b_j = -b_{j-1} - p \frac{\partial}{\partial p} b_{j-1} - q \frac{\partial}{\partial p} b_{j-2},$$

$$\frac{\partial}{\partial q} b_j = -b_{j-2} - p \frac{\partial}{\partial q} b_{j-1} - q \frac{\partial}{\partial q} b_{j-2}.$$

Теперь определим величины c_j при помощи рекурсионных равенств

$$c_j = b_j - pc_{j-1} - qc_{j-2} \quad (j=0, 1, \dots, n-1), \quad (1.12.8)$$

$$c_{-1} = c_{-2} = 0,$$

откуда они могут быть найдены последовательно.

Если сравнить эти равенства с соотношениями для частных производных $\frac{\partial}{\partial p} b_j$ и $\frac{\partial}{\partial q} b_j$, записанными выше, можно сказать, что

$$\frac{\partial}{\partial p} b_j = -c_{j-1}, \quad \frac{\partial}{\partial q} b_j = -c_{j-2} \quad (j=0, 1, 2, \dots, n), \quad (1.12.9)$$

$$c_{-1} = c_{-2} = 0.$$

Таким образом, равенства (1.12.8) дают правило рекурсионного нахождения частных производных $\frac{\partial}{\partial p} b_j$ и $\frac{\partial}{\partial q} b_j$ по значениям p, q и ранее найденным значениям $b_j(p, q)$.

Сравним равенства (1.12.8) с (1.12.7). Числа c_j получаются из b_j по таким же правилам, как b_j получаются из a_j . Поэтому c_j могут быть найдены при помощи деления $L(x) = xQ(x) + b_{n-1} = b_0x^{n-1} + b_1x^{n-2} + \dots + b_{n-1}$ на $x^2 + px + q$ с условием особого представления остатка:

$$L(x) = (x^2 + px + q)(c_0x^{n-3} + c_1x^{n-4} + \dots + c_{n-3}) + c_{n-2}(x + p) + c_{n-1}.$$

Литература

1. Березин И. С., Жидков Н. П. Методы вычислений. М., 1966.
2. Загускин В. Л. Справочник по численным методам решения уравнений. М., 1960.
3. Канторович Л. В., Акилов Г. П. Функциональный анализ в нормированных пространствах, гл. XVII, XVIII. М., 1959.
4. Коллатц Л. Функциональный анализ и вычислительная математика. М., 1969.
5. Ланс Дж. Н. Численные методы для быстродействующих вычислительных машин. М., 1962.
6. Мысовских И. П. Лекции по методам вычислений. М., 1962.
7. Островский А. Решение уравнений и систем уравнений. М., 1963.
8. Хаусхолдер А. С. Основы численного анализа. М., 1956.

Глава 2

РЕШЕНИЕ СИСТЕМ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ

В этой главе будут рассмотрены простейшие и уже давно испытанные методы решения систем линейных алгебраических уравнений. К решению таких систем приводятся многие прикладные задачи, а также значительная часть задач численного анализа. Наряду с проблемой решения неоднородной системы линейных алгебраических уравнений здесь будет изучена и тесно связанная с ней проблема обращения матрицы, а также задача вычисления определителя матрицы.

Методы решения систем линейных алгебраических уравнений, на изучении которых мы здесь остановимся, можно разделить на две большие группы: точные и итерационные. Под точными мы будем подразумевать такие методы, которые позволяют получить точные значения неизвестных в результате выполнения конечного числа арифметических операций. Хорошо известное из курса линейной алгебры правило Крамера может служить примером такого метода. Правда, правило Крамера в практике решения систем линейных алгебраических уравнений обычно не применяется, так как оно требует выполнения очень большого количества арифметических операций и записей. Здесь мы будем рассматривать более экономичные точные методы, значительная часть которых основана на знакомой еще по школьному курсу математики идее последовательного исключения неизвестных из уравнений системы.

Итерационные же методы решения систем линейных алгебраических уравнений характеризуются тем, что точное решение системы они могут, вообще говоря, давать лишь как предел некоторой бесконечной последовательности векторов. Исходное приближение к решению при этом разывается каким-либо другим способом или задается произвольно. При выполнении определенных требований можно получить достаточно быстро сходящийся к решению итерационный процесс.

Прежде чем приступить к рассмотрению конкретных методов решения, приведем здесь некоторые сведения из линейной алгебры, которые будут существенно использоваться в дальнейшем, особенно при изучении итерационных процессов.

§ 2.1. НЕКОТОРЫЕ СВЕДЕНИЯ ИЗ ЛИНЕЙНОЙ АЛГЕБРЫ

2.1.1. Сходимость последовательностей векторов и матриц

При описании итерационных методов решения систем линейных алгебраических уравнений нам в первую очередь понадобится понятие предела для объектов линейной алгебры. Так как численные задачи линейной алгебры обычно формулируются в терминах матриц, то мы определим понятие предела для однострочных матриц, которые будем отождествлять с векторами арифметического пространства, и для квадратных матриц. При этом для удобства записи мы будем иногда вектор-столбец представлять в виде транспонированной однострочной матрицы.

Пусть дана последовательность векторов

$$\bar{x}^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})' \quad (k=0, 1, 2, \dots).$$

Если существуют n конечных пределов

$$x_i = \lim_{k \rightarrow \infty} x_i^{(k)} \quad (i=1, 2, \dots, n),$$

то вектор $\bar{x} = (x_1, x_2, \dots, x_n)'$ называют пределом последовательности $\bar{x}^{(k)}$ ($k=0, 1, 2, \dots$), а саму эту последовательность называют сходящейся к вектору \bar{x} .

Точно так же, если имеется последовательность квадратных матриц $A^{(k)} = (a_{ij}^{(k)})$ ($i, j=1, 2, \dots, n$; $k=0, 1, \dots$), то пределом этой последовательности называют матрицу A с элементами $a_{ij} = \lim_{k \rightarrow \infty} a_{ij}^{(k)}$, если, конечно, все эти n^2 пределов существуют.

В соответствии с таким определением предела бесконечный ряд матриц называют сходящимся, если существует предел последовательности его частных сумм. Предел этот и называют суммой данного ряда. Очевидно, что ряд матриц будет сходящимся тогда и только тогда, когда будут сходиться все ряды из одноименных элементов, при этом суммы этих рядов будут являться элементами суммы данного ряда матриц.

2.1.2. Нормы векторов и матриц

Введем сначала понятие нормы вектора, обобщающее известное понятие длины вектора. Назовем *нормой вектора* \bar{x} сопоставляемое этому вектору вещественное число $\|\bar{x}\|$, удовлетворяющее следующим требованиям:

- 1) $\|\bar{x}\| > 0$ при $\bar{x} \neq \bar{0}$ и $\|\bar{0}\| = 0$;
- 2) $\|c\bar{x}\| = |c| \cdot \|\bar{x}\|$ при любом числовом множителе c ;
- 3) $\|\bar{x} + \bar{y}\| \leq \|\bar{x}\| + \|\bar{y}\|$.

Из такого определения нормы вектора непосредственно следует, что

$$\|\bar{x}-\bar{y}\| \geq \left| \|\bar{x}\| - \|\bar{y}\| \right|.$$

В самом деле,

$$\|\bar{x}\| = \|\bar{x}-\bar{y}+\bar{y}\| \leq \|\bar{x}-\bar{y}\| + \|\bar{y}\|,$$

т. е.

$$\|\bar{x}-\bar{y}\| \geq \|\bar{x}\| - \|\bar{y}\|.$$

И аналогично

$$\|\bar{x}-\bar{y}\| = \|\bar{y}-\bar{x}\| \geq \|\bar{y}\| - \|\bar{x}\|,$$

или

$$\|\bar{x}-\bar{y}\| \geq -(\|\bar{x}\| - \|\bar{y}\|).$$

Следовательно,

$$\|\bar{x}-\bar{y}\| \geq \left| \|\bar{x}\| - \|\bar{y}\| \right|.$$

Вводить норму вектора можно различными способами, только бы выполнялись условия (1)–(3) данного выше определения нормы. Приведем примеры наиболее распространенных способов задания нормы вектора $\bar{x} = (x_1, x_2, \dots, x_n)'$.

1. Первая (кубическая) норма

$$\|\bar{x}\|_I = \max_{1 \leq i \leq n} |x_i|.$$

Введенную так норму обычно называют кубической в связи с тем, что множество векторов вещественного пространства, для которых $\|\bar{x}\|_I \leq 1$, заполняет единичный куб

$$-1 \leq x_i \leq 1 \quad (i=1, 2, \dots, n).$$

2. Вторая (октаэдрическая) норма

$$\|\bar{x}\|_{II} = \sum_{i=1}^n |x_i|.$$

Множество вещественных векторов, для которых $\|\bar{x}\|_{II} \leq 1$, заполняет n -мерный аналог октаэдра.

3. Третья (сферическая или евклидова) норма

$$\|\bar{x}\|_{\text{III}} = \sqrt{(\bar{x}, \bar{x})} = \sqrt{\sum_{i=1}^n |x_i|^2} = |\bar{x}|.$$

Третья норма вектора есть не что иное, как длина вектора. Совокупность векторов, для которых $\|\bar{x}\|_{\text{III}} \leq 1$, заполняет шар единичного радиуса.

Для этих трех норм выполнимость всех условий, данных в определении нормы вектора, легко проверяется.

Дадим теперь иное определение сходимости последовательности векторов, основанное на введенном понятии нормы, а именно, будем говорить, что $\bar{x}^{(k)} \xrightarrow[k \rightarrow \infty]{} \bar{x}$, если $\|\bar{x} - \bar{x}^{(k)}\| \xrightarrow[k \rightarrow \infty]{} 0$. Эквивалентность такого определения сходимости по норме прежнему определению сходимости в координатах основана на следующем утверждении.

Теорема 1. Для того чтобы $\bar{x}^{(k)} \xrightarrow[k \rightarrow \infty]{} \bar{x}$, необходимо и достаточно, чтобы $\|\bar{x} - \bar{x}^{(k)}\| \xrightarrow[k \rightarrow \infty]{} 0$.

Доказательство. Проверим сначала необходимость высказанного условия. Пусть $\bar{x}^{(k)} \xrightarrow[k \rightarrow \infty]{} \bar{x}$, т. е. $\lim_{k \rightarrow \infty} x_i^{(k)} = x_i$ для всех $i = 1, 2, \dots, n$. Введя векторы $\bar{e}_1 = (1, 0, 0, \dots, 0)'$, $\bar{e}_2 = (0, 1, 0, \dots, 0)'$, ..., $\bar{e}_n = (0, 0, 0, \dots, 1)'$, можем записать:

$$\bar{x} - \bar{x}^{(k)} = \sum_{i=1}^n (x_i - x_i^{(k)}) \bar{e}_i \quad (k=0, 1, 2, \dots).$$

Если через N обозначить $\max_{1 \leq i \leq n} \|\bar{e}_i\|$, то из последних равенств следует, что

$$\|\bar{x} - \bar{x}^{(k)}\| \leq N \sum_{i=1}^n |x_i - x_i^{(k)}| \quad (k=0, 1, 2, \dots).$$

Поэтому

$$\|\bar{x} - \bar{x}^{(k)}\| \xrightarrow[k \rightarrow \infty]{} 0.$$

Теперь проверим достаточность высказанного в теореме условия. Пусть $\lim_{k \rightarrow \infty} \|\bar{x} - \bar{x}^{(k)}\| = 0$. Тогда, так как

$$\|\bar{x}^{(k)}\| = \|\bar{x} + (\bar{x}^{(k)} - \bar{x})\| \leq \|\bar{x}\| + \|\bar{x} - \bar{x}^{(k)}\|,$$

*) Скалярное произведение векторов вводится по формуле $(\bar{x}, \bar{y}) = \sum_{i=1}^n x_i \bar{y}_i$, где через \bar{y}_i обозначено число, комплексно сопряженное с координатой y_i вектора \bar{y} .

то $\|\bar{x}^{(k)}\|$ будет ограниченной при любом $k=0, 1, 2, \dots$, т. е. $\|\bar{x}^{(k)}\| \leq M$ ($k=0, 1, 2, \dots$). Покажем, что будет ограниченной при любом $k=0, 1, 2, \dots$ также и величина $c_k = |x_1^{(k)}| + |x_2^{(k)}| + \dots + |x_n^{(k)}|$, т. е. что $c_k \leq L$ ($k=0, 1, 2, \dots$). Предположим противное, т. е. допустим, что существует такая последовательность k_1, k_2, \dots индексов, что $c_{k_m} \xrightarrow{m \rightarrow \infty} \infty$. Для простоты записи будем считать, что $c_k \xrightarrow{k \rightarrow \infty} \infty$. По данной последовательности векторов $\bar{x}^{(k)}$ ($k=0, 1, 2, \dots$) строим последовательность векторов

$$\bar{y}^{(k)} = \frac{\bar{x}^{(k)}}{c_k} = (y_1^{(k)}, y_2^{(k)}, \dots, y_n^{(k)})' \quad (k=0, 1, 2, \dots).$$

Так как $y_i^{(k)} = \frac{x_i^{(k)}}{c_k}$ ($i=1, 2, \dots, n$), то

$$|y_1^{(k)}| + |y_2^{(k)}| + \dots + |y_n^{(k)}| = 1 \quad (k=0, 1, 2, \dots),$$

и координаты векторов $\bar{y}^{(k)}$ ограничены в совокупности. Поэтому мы можем выбрать такую последовательность индексов, что будут существовать конечные пределы

$$\lim_{k \rightarrow \infty} y_i^{(k)} = y_i \quad (i=1, 2, \dots, n).$$

Так как $|y_1| + |y_2| + \dots + |y_n| = 1$, то предельный вектор

$$\bar{y} = (y_1, y_2, \dots, y_n)'$$

отличен от нулевого вектора. С другой стороны, так как

$$\|\bar{y}\| = \|\bar{y}^{(k)} + (\bar{y} - \bar{y}^{(k)})\| \leq \|\bar{y}^{(k)}\| + \|\bar{y} - \bar{y}^{(k)}\| = \frac{\|\bar{x}^{(k)}\|}{c_k} + \|\bar{y} - \bar{y}^{(k)}\| \quad (k=0, 1, 2, \dots),$$

при этом $\|\bar{x}^{(k)}\| \leq M$ по доказанному, $c_k \xrightarrow{k \rightarrow \infty} \infty$ по предположению, а $\|\bar{y} - \bar{y}^{(k)}\| \xrightarrow{k \rightarrow \infty} 0$ в силу уже доказанной необходимости условий теоремы, то $\|\bar{y}\| = 0$, т. е. $\bar{y} = \bar{0}$. Полученное противоречие доказывает, что $c_k \leq L$ ($k=0, 1, \dots$), т. е. что координаты векторов $\bar{x}^{(k)}$ ограничены в совокупности. Это позволяет выбрать такую последовательность индексов, что будут существовать конечные пределы $\xi_i = \lim_{k \rightarrow \infty} x_i^{(k)}$ ($i=1, 2, \dots, n$). Покажем, что предельный вектор $\bar{\xi} = (\xi_1, \xi_2, \dots, \xi_n)'$ совпадает с вектором $\bar{x} = (x_1, x_2, \dots, x_n)'$. В самом деле, так как $\|\bar{x} - \bar{x}^{(k)}\| \xrightarrow{k \rightarrow \infty} 0$ по условию, $\|\bar{\xi} - \bar{x}^{(k)}\| \xrightarrow{k \rightarrow \infty} 0$ в силу доказанной необходимости условий теоремы, а $\|\bar{x} - \bar{\xi}\| = \|(\bar{x} - \bar{x}^{(k)}) + (\bar{x}^{(k)} - \bar{\xi})\| \leq \|\bar{x} - \bar{x}^{(k)}\| + \|\bar{x}^{(k)} - \bar{\xi}\|$ для любого $k=0, 1, 2, \dots$, то $\|\bar{x} - \bar{\xi}\| = 0$, т. е. $\bar{\xi} = \bar{x}$.
Теорема доказана.

Заметим также, что из условия $\bar{x}^{(k)} \xrightarrow{k \rightarrow \infty} \bar{x}$ следует, что $\|\bar{x}^{(k)}\| \xrightarrow{k \rightarrow \infty} \|\bar{x}\|$.

Справедливость этого утверждения становится очевидной, если учесть неравенство

$$\left| \|\bar{x}\| - \|\bar{x}^{(k)}\| \right| \leq \|\bar{x} - \bar{x}^{(k)}\|.$$

Рассмотрим далее понятие нормы матрицы. *Нормой квадратной матрицы* A назовем сопоставляемое ей вещественное число $\|A\|$, удовлетворяющее следующим условиям:

- 1) $\|A\| > 0$ при $A \neq O$ и $\|O\| = 0$,
- 2) $\|cA\| = |c| \cdot \|A\|$ при любом числовом множителе c ,
- 3) $\|A+B\| \leq \|A\| + \|B\|$,
- 4) $\|AB\| \leq \|A\| \cdot \|B\|$.

Аналогично случаю вектора можно показать, что из такого определения нормы матрицы следует неравенство

$$\|A-B\| \geq \left| \|A\| - \|B\| \right|.$$

Так же как и в случае векторов, для последовательности матриц можно дать иное определение сходимости, доказав, что условие $\|A - A^{(k)}\| \xrightarrow[k \rightarrow \infty]{} 0$ является необходимым и достаточным условием того, что $A^{(k)} \xrightarrow[k \rightarrow \infty]{} A$.

Из неравенства $\left| \|A\| - \|B\| \right| \leq \|A-B\|$ теперь уже непосредственно следует, что, если $A^{(k)} \xrightarrow[k \rightarrow \infty]{} A$, то $\|A^{(k)}\| \xrightarrow[k \rightarrow \infty]{} \|A\|$.

Норма матрицы также может быть задана многими способами. Однако поскольку в большинстве задач линейной алгебры обычно в рассуждении одновременно участвуют как матрицы, так и векторы, то норму матрицы целесообразно вводить так, чтобы она была разумным образом связана с данной нормой вектора. Будем говорить, что норма матриц *согласована* с данной нормой векторов, если для любой квадратной матрицы A и для любого вектора \bar{x} , размерность которого равна порядку матрицы, выполняется неравенство

$$\|A\bar{x}\| \leq \|A\| \cdot \|\bar{x}\|.$$

Среди всех норм матриц, согласованных с данной нормой векторов, выберем наименьшую. Для этих целей за норму матрицы A примем максимум норм векторов $A\bar{x}$ в предположении, что вектор \bar{x} пробегает множество всех векторов, норма которых равна единице:

$$\|A\| = \max_{\|\bar{x}\|=1} \|A\bar{x}\|.$$

Для каждой матрицы A в силу непрерывности нормы этот максимум достигается, т. е. всегда найдется вектор $\bar{x}^{(0)}$ такой, что $\|\bar{x}^{(0)}\|=1$ и $\|A\bar{x}^{(0)}\| = \|A\|$. Введенную так норму матриц будем называть *подчиненной* данной нормой векторов.

Проверим, что норма матриц, подчиненная данной норме векторов, во-первых, удовлетворяет условиям (1)—(4) определения нормы матриц, во-вторых, согласована с этой нормой векторов и, в-третьих, не больше всякой нормы, согласованной с той же нормой векторов.

Начнем с проверки условия (1), данного в определении нормы матриц. Пусть $A \neq 0$. Тогда найдется вектор \bar{y} такой, что $A\bar{y} \neq \bar{0}$. По вектору \bar{y} построим вектор $\bar{x} = \frac{\bar{y}}{\|\bar{y}\|}$, для которого уже выполняется требование $\|\bar{x}\|=1$. Так как $A\bar{x} \neq \bar{0}$, то $\|A\bar{x}\| > 0$, значит и

$$\|A\| = \max_{\|\bar{x}\|=1} \|A\bar{x}\| > 0.$$

Если же $A = O$, то $\|A\| = \max_{\|\bar{x}\|=1} \|O\bar{x}\| = 0$.

Справедливость условия (2) проверяется непосредственно:

$$\|cA\| = \max_{\|\bar{x}\|=1} \|cA\bar{x}\| = \max_{\|\bar{x}\|=1} |c| \cdot \|A\bar{x}\| = |c| \max_{\|\bar{x}\|=1} \|A\bar{x}\| = |c| \cdot \|A\|.$$

Проверим условие (3). Как мы уже отмечали, для каждой матрицы $A+B$ всегда найдется вектор $\bar{x}^{(0)}$ такой, что $\|\bar{x}^{(0)}\|=1$ и

$$\|A+B\| = \max_{\|\bar{x}\|=1} \|(A+B)\bar{x}\| = \|(A+B)\bar{x}^{(0)}\|.$$

Тогда

$$\begin{aligned} \|A+B\| &= \|A\bar{x}^{(0)} + B\bar{x}^{(0)}\| \leq \|A\bar{x}^{(0)}\| + \|B\bar{x}^{(0)}\| \leq \\ &\leq \max_{\|\bar{x}\|=1} \|A\bar{x}\| + \max_{\|\bar{x}\|=1} \|B\bar{x}\| = \|A\| + \|B\|. \end{aligned}$$

Прежде чем проверить условие (4), установим выполнимость требования согласования

$$\|A\bar{x}\| \leq \|A\| \cdot \|\bar{x}\|.$$

Если $\bar{x} = \bar{0}$, то справедливость этого неравенства очевидна.

Пусть $\bar{x} \neq \bar{0}$. Тогда рассмотрим вектор $\bar{y}^{(0)} = \frac{\bar{x}}{\|\bar{x}\|}$. Так как $\|\bar{y}^{(0)}\|=1$, то

$$\|A\bar{x}\| = \|A(\|\bar{x}\|\bar{y}^{(0)})\| = \|\bar{x}\| \cdot \|A\bar{y}^{(0)}\| \leq \|\bar{x}\| \max_{\|\bar{y}\|=1} \|A\bar{y}\| = \|A\| \cdot \|\bar{x}\|.$$

Проверим, наконец, условие (4). Как и прежде, для матрицы AB найдем вектор $\bar{x}^{(0)}$ такой, что

$$\|\bar{x}^{(0)}\| = 1 \text{ и } \|AB\bar{x}^{(0)}\| = \|AB\|.$$

Тогда

$$\|AB\| = \|A(B\bar{x}^{(0)})\| \leq \|A\| \cdot \|B\bar{x}^{(0)}\| \leq \|A\| \cdot \|B\| \cdot \|\bar{x}^{(0)}\| = \|A\| \cdot \|B\|.$$

Осталось проверить лишь утверждение о том, что норма матриц, подчиненная данной норме векторов, не больше любой нормы, согласованной с той же нормой векторов. Действительно, пусть $\|A\|$ есть норма матрицы A , подчиненная данной норме векторов, а $\|A\|_c$ есть любая норма матрицы A , согласованная с той же нормой векторов. Тогда, как мы уже знаем, для матрицы A найдется вектор $\bar{x}^{(0)}$ такой, что

$$\|\bar{x}^{(0)}\| = 1 \text{ и } \|A\bar{x}^{(0)}\| = \|A\|.$$

Но

$$\|A\bar{x}^{(0)}\| \leq \|A\|_c \cdot \|\bar{x}^{(0)}\| = \|A\|_c$$

и, значит,

$$\|A\| \leq \|A\|_c.$$

Для любой нормы матриц, подчиненной норме векторов, $\|E\| = 1$. Здесь и всюду в гл. 2 и 3 через E обозначена единичная матрица.

Рассмотрим два примера задания нормы матрицы $A = (a_{ij})$. Положим

$$M(A) = n \max_{1 \leq i, j \leq n} |a_{ij}| \text{ и } N(A) = \sqrt{\text{Sp } A^*A},$$

где A^* есть матрица, сопряженная с матрицей A , т. е. комплексно сопряженная с транспонированной матрицей A' , а след матрицы A^*A равен сумме диагональных элементов этой матрицы, т. е.

$$\text{Sp } A^*A = \sum_{j=1}^n (A^*A)_{jj} = \sum_{i,j=1}^n |a_{ij}|^2.$$

Для $M(A)$ и $N(A)$ легко проверяется выполнимость всех четырех условий определения нормы матрицы. Нетрудно также установить, что норма $M(A)$ согласована с кубической, октаэдрической и сферической нормами вектора, а норма $N(A)$ согласована со сферической нормой вектора. Однако ни норма $M(A)$, ни норма $N(A)$ не подчинены ни одной из норм векторов, так как $M(E) = n$, а $N(E) = \sqrt{n}$.

Укажем представления норм матриц, подчиненных введенным ранее трем нормам векторов.

Рассмотрим сначала первую (кубическую) норму векторов

$$\|\bar{x}\|_I = \max_{1 \leq i \leq n} |x_i|.$$

Оказывается, подчиненная ей норма матрицы A такова:

$$\|A\|_I = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

Действительно, для любого вектора \bar{x} единичной кубической нормы

$$\|A\bar{x}\|_I = \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij}x_j \right| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \cdot |x_j| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$$

и, следовательно,

$$\|A\|_I = \max_{\|\bar{x}\|_I=1} \|A\bar{x}\|_I \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

С другой стороны, возьмем вектор $\bar{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})'$, где

$$x_j^{(0)} = \begin{cases} \frac{|a_{kj}|}{a_{kj}}, & \text{если } a_{kj} \neq 0; \\ 1, & \text{если } a_{kj} = 0 \end{cases} \quad (j=1, 2, \dots, n),$$

а k есть номер строки матрицы A , на которой достигается

$$\max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

Очевидно, что $\|\bar{x}^{(0)}\|_I = 1$. Кроме того,

$$\left| \sum_{j=1}^n a_{ij}x_j^{(0)} \right| \leq \sum_{j=1}^n |a_{ij}| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| = \sum_{j=1}^n |a_{kj}|$$

при $i \neq k$ и

$$\left| \sum_{j=1}^n a_{kj}x_k^{(0)} \right| = \sum_{j=1}^n |a_{kj}|.$$

Тогда и

$$\|A\bar{x}^{(0)}\|_I = \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij}x_j^{(0)} \right| = \sum_{j=1}^n |a_{kj}| = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

Следовательно,

$$\|A\|_I = \max_{\|\bar{x}\|_I=1} \|A\bar{x}\|_I \geq \|A\bar{x}^{(0)}\|_I = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

Сопоставление неравенств

$$\|A\|_I \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad \text{и} \quad \|A\|_I \geq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$$

убеждает нас в справедливости высказанного утверждения.

Укажем, далее, норму матриц, подчиненную второй (октаэдрической) норме векторов

$$\|\bar{x}\|_{II} = \sum_{i=1}^n |x_i|.$$

Такой нормой, оказывается, будет

$$\|A\|_{II} = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|.$$

В самом деле, если $\|\bar{x}\|_{II} = 1$, то

$$\begin{aligned} \|A\bar{x}\|_{II} &= \sum_{i=1}^n \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \sum_{i=1}^n \sum_{j=1}^n |a_{ij}| \cdot |x_j| = \\ &= \sum_{j=1}^n |x_j| \left(\sum_{i=1}^n |a_{ij}| \right) \leq \left(\max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| \right) \cdot \sum_{j=1}^n |x_j| = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|. \end{aligned}$$

Пусть $\sum_{i=1}^n |a_{ij}|$ достигает наибольшего значения для столбца с номером k . Тогда вектор

$$\bar{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})',$$

где $x_i^{(0)} = 0$ при $i \neq k$ и $x_k^{(0)} = 1$, имеет октаэдрическую норму, равную единице, и

$$\|A\bar{x}^{(0)}\|_{II} = \sum_{i=1}^n \left| \sum_{j=1}^n a_{ij} x_j^{(0)} \right| = \sum_{i=1}^n |a_{ik}| = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|.$$

Следовательно,

$$\max_{\|\bar{x}\|_{II}=1} \|A\bar{x}\|_{II} = \|A\bar{x}^{(0)}\|_{II} = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|,$$

т. е.

$$\|A\|_{II} = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|.$$

Рассмотрим, наконец, третью (сферическую) норму векторов

$$\|\bar{x}\|_{III} = \sqrt[n]{|\bar{x}|} = \sqrt[n]{(\bar{x}, \bar{x})} = \sqrt[n]{\sum_{i=1}^n |x_i|^2}.$$

Оказывается, что подчиненной ей нормой матриц будет

$$\|A\|_{\text{III}} = \sqrt[n]{\Lambda_1},$$

где Λ_1 есть наибольшее собственное значение матрицы A^*A (напомним, что число β называется собственным значением матрицы B , если существует ненулевой вектор \bar{x} , для которого $B\bar{x} = \beta\bar{x}$; вектор \bar{x} называется в этом случае собственным вектором матрицы B , отвечающим данному собственному значению β).

Прежде чем приступить к доказательству высказанного утверждения, необходимо проверить, что Λ_1 будет неотрицательным числом.

Вещественность Λ_1 уже следует из того, что матрица A^*A есть эрмитова (матрица B называется эрмитовой, если $B^* = B$). В самом деле,

$$(A^*A)^* = (A)^*(A^*)^* = A^*A.$$

Эрмитова же матрица, как известно, имеет только вещественные собственные значения и обладает полной системой попарно ортогональных собственных векторов. Проверим, что Λ_1 будет неотрицательным числом. Действительно, для любого собственного значения Λ матрицы A^*A существует такой вектор $\bar{x} \neq 0$, что $A^*A\bar{x} = \Lambda\bar{x}$. Умножая обе части этого равенства скалярно на вектор \bar{x} и учитывая свойства скалярного произведения векторов, получим

$$(A^*A\bar{x}, \bar{x}) = (\Lambda\bar{x}, \bar{x}). \quad (2.1.1)$$

При этом

$$(A^*A\bar{x}, \bar{x}) = (A\bar{x}, A\bar{x}) = \|A\bar{x}\|_{\text{III}}^2$$

и

$$(\Lambda\bar{x}, \bar{x}) = \Lambda(\bar{x}, \bar{x}) = \Lambda \|\bar{x}\|_{\text{III}}^2,$$

так что равенство (2.1.1) можно записать в виде

$$\|A\bar{x}\|_{\text{III}}^2 = \Lambda \|\bar{x}\|_{\text{III}}^2,$$

откуда и следует, что $\Lambda \geq 0$.

Докажем, наконец, что $\|A\|_{\text{III}} = \sqrt[n]{\Lambda_1}$. В самом деле, пусть

$$\Lambda_1 \geq \Lambda_2 \geq \dots \geq \Lambda_n$$

есть собственные значения матрицы A^*A , а

$$\bar{x}^{(1)}, \bar{x}^{(2)}, \dots, \bar{x}^{(n)}$$

есть соответствующие им векторы полной системы собственных векторов этой матрицы. Будем считать эти векторы ортонормированными. Рассмотрим любой вектор \bar{x} единичной евклидовой нормы и разложим его по собственным векторам $\bar{x}^{(1)}, \bar{x}^{(2)}, \dots, \bar{x}^{(n)}$:

$$\bar{x} = \sum_{i=1}^n \alpha_i \bar{x}^{(i)}.$$

Так как $\|\bar{x}\|_{III} = 1$, то

$$\|\bar{x}\|_{III}^2 = (\bar{x}, \bar{x}) = \sum_{i=1}^n |\alpha_i|^2 = 1.$$

Тогда

$$\begin{aligned} \|A\bar{x}\|_{III}^2 &= (A\bar{x}, A\bar{x}) = (\bar{x}, A^*A\bar{x}) = \left(\sum_{i=1}^n \alpha_i \bar{x}^{(i)}, \sum_{i=1}^n \alpha_i \Lambda_i \bar{x}^{(i)} \right) = \\ &= \sum_{i=1}^n |\alpha_i|^2 \Lambda_i \leq \Lambda_1 \sum_{i=1}^n |\alpha_i|^2 = \Lambda_1. \end{aligned}$$

Поэтому

$$\|A\|_{III} = \max_{\|\bar{x}\|_{III}=1} \|A\bar{x}\|_{III} \leq \sqrt{\Lambda_1}.$$

С другой стороны, если взять в качестве вектора \bar{x} вектор $\bar{x}^{(1)}$, то

$$\|A\bar{x}^{(1)}\|_{III}^2 = (\bar{x}^{(1)}, A^*A\bar{x}^{(1)}) = (\bar{x}^{(1)}, \Lambda_1 \bar{x}^{(1)}) = \Lambda_1,$$

и

$$\|A\|_{III} = \max_{\|\bar{x}\|_{III}=1} \|A\bar{x}\|_{III} \geq \sqrt{\Lambda_1}.$$

Значит,

$$\|A\|_{III} = \sqrt{\Lambda_1}.$$

В частном случае, когда матрица A есть эрмитова матрица и, стало быть, $A^*A = A^2$, будет $\Lambda_1 = \lambda_1^2$, где λ_1 есть наибольшее по модулю собственное значение матрицы A , и

$$\|A\|_{III} = |\lambda_1|.$$

Сравнение введенных норм матриц приводит к следующим неравенствам, справедливость которых здесь не проверяется:

$$\frac{1}{n} M(A) \leq \|A\|_i \leq M(A) \quad (i=I, II, III),$$

$$\frac{1}{n} M(A) \leq N(A) \leq M(A),$$

$$\frac{1}{\sqrt{n}} N(A) \leq \|A\|_{III} \leq N(A),^*)$$

$$\frac{1}{\sqrt{n}} N(A) \leq \|A\|_j \leq \sqrt{n} N(A) \quad (j=I, II),$$

*) Последнее неравенство было доказано при получении формулы (1.10.18).

$$\frac{1}{\sqrt[n]{n}} \|A\|_{III} \leq \|A\|_k \leq \sqrt[n]{n} \|A\|_{III} \quad (k=I, II),$$

$$\frac{1}{n} \|A\|_I \leq \|A\|_{II} \leq n \|A\|_I.$$

Подобные же неравенства можно записать и для ранее введенных норм векторов:

$$\|\bar{x}\|_I \leq \|\bar{x}\|_{II} \leq n \|\bar{x}\|_I,$$

$$\|\bar{x}\|_I \leq \|\bar{x}\|_{III} \leq \sqrt[n]{n} \|\bar{x}\|_I,$$

$$\frac{1}{\sqrt[n]{n}} \|\bar{x}\|_{II} \leq \|\bar{x}\|_{III} \leq \|\bar{x}\|_{II}.$$

2.1.3. Сходимость матричной геометрической прогрессии

Рассмотрим матричный ряд

$$E + A + A^2 + \dots + A^m + \dots$$

Встает вопрос, при каких условиях эта матричная геометрическая прогрессия сходится и, если она сходится, чему равна ее сумма. Если бы мы имели дело с обычной числовой геометрической прогрессией

$$1 + a + a^2 + \dots + a^m + \dots,$$

то необходимым и достаточным условием ее сходимости было бы условие $a^m \xrightarrow{m \rightarrow \infty} 0$, при этом ее сумма была бы равна

$$(1-a)^{-1}.$$

Оказывается, подобные же результаты имеют место и в случае матричной геометрической прогрессии. Прежде чем высказать их, сформулируем некоторые утверждения предварительного характера.

Лемма 1. Для того чтобы $A^m \xrightarrow{m \rightarrow \infty} 0$, необходимо и достаточно, чтобы все собственные значения матрицы A были по модулю меньше единицы.

Доказательство. Известно, что с помощью преобразования подобия, которое не меняет собственных значений матрицы, исходная матрица A всегда может быть приведена к канонической форме Жордана

$$I = C^{-1}AC.$$

Здесь C — некоторая матрица, а I — квазидиагональная матрица

$[I_{\tau_1}(\lambda_1), I_{\tau_2}(\lambda_2), \dots, I_{\tau_r}(\lambda_r)]$, где r — число канонических ящиков Жордана

$$I_{\tau}(\lambda) = \begin{bmatrix} \lambda & 0 & 0 & \dots & 0 & 0 \\ 1 & \lambda & 0 & \dots & 0 & 0 \\ 0 & 1 & \lambda & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & \lambda \end{bmatrix}.$$

r равно числу линейно независимых собственных векторов матрицы A , при этом $\sum_{i=1}^r \tau_i = n$, где τ_i — порядок i -го ящика Жордана $I_{\tau_i}(\lambda_i)$, а n — порядок исходной матрицы A . Тогда

$$A = C I C^{-1},$$

$$A^m = C I C^{-1} \cdot C I C^{-1} \cdot \dots \cdot C I C^{-1} = C I^m C^{-1}.$$

Поэтому матрицы A^m и I^m при $m \rightarrow \infty$ одновременно стремятся или не стремятся к нулевой матрице.

Так как

$$I^m = [I_{\tau_1}^m(\lambda_1), I_{\tau_2}^m(\lambda_2), \dots, I_{\tau_r}^m(\lambda_r)],$$

то для выяснения условий сходимости $A^m \xrightarrow{m \rightarrow \infty} 0$ достаточно установить

лишь условия сходимости $I_{\tau}^m(\lambda) \xrightarrow{m \rightarrow \infty} 0$. Непосредственной проверкой лег-

ко убедиться в справедливости следующего равенства:

$$I_{\tau}^m(\lambda) = \begin{bmatrix} \lambda^m & 0 & 0 & \dots & 0 & 0 \\ \frac{(\lambda^m)'}{1!} & \lambda^m & 0 & \dots & 0 & 0 \\ \frac{(\lambda^m)''}{2!} & \frac{(\lambda^m)'}{1!} & \lambda^m & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \frac{(\lambda^m)^{(\tau-1)}}{(\tau-1)!} & \frac{(\lambda^m)^{(\tau-2)}}{(\tau-2)!} & \frac{(\lambda^m)^{(\tau-3)}}{(\tau-3)!} & \dots & \frac{(\lambda^m)'}{1!} & \lambda^m \end{bmatrix},$$

где для удобства записи условно введена операция дифференцирования по λ . Диагональными элементами матрицы $I_{\tau}^m(\lambda)$ являются λ^m , поэтому для сходимости при $m \rightarrow \infty$ последовательности матриц $I_{\tau}^m(\lambda)$ к нулевой матрице необходимо, чтобы выполнялось условие $|\lambda| < 1$. Но выполнения этого условия и достаточно для сходимости $I_{\tau}^m(\lambda) \xrightarrow{m \rightarrow \infty} 0$, ибо тогда

$$\frac{(\lambda^m)^{(i)}}{i!} \xrightarrow{m \rightarrow \infty} 0$$

для любого $i=0, 1, \dots, \tau-1$.

Признак сходимости $A^m \xrightarrow{m \rightarrow \infty} O$, сформулированный в только что доказанной лемме 1, неудобен для проверки, так как требует наличия достаточно точной информации о собственных значениях матрицы A . Более удобным в этом отношении является следующий признак.

Лемма 2. Для того чтобы $A^m \xrightarrow{m \rightarrow \infty} O$, достаточно, чтобы хоть одна из норм матрицы A была меньше единицы.

Доказательство. Как нам уже известно, для того чтобы установить, что $A^m \xrightarrow{m \rightarrow \infty} O$, достаточно проверить, что $\|O - A^m\| \xrightarrow{m \rightarrow \infty} 0$ хотя бы

для одной из норм матрицы A . Но

$$\|O - A^m\| = \|A^m\| = \|A^{m-1} \cdot A\| \leq \|A^{m-1}\| \cdot \|A\| \leq \dots \leq \|A\|^m.$$

Поэтому, если какая-либо $\|A\| < 1$, то $\|A^m\| \xrightarrow{m \rightarrow \infty} 0$, т. е. $A^m \xrightarrow{m \rightarrow \infty} O$, что и требовалось доказать.

Опираясь на леммы 1 и 2, можно доказать следующее утверждение о сравнительной величине нормы матрицы и ее собственных значений, которое мы также будем использовать в дальнейшем.

Лемма 3. Модуль каждого собственного значения матрицы не превосходит любой из ее норм.

Доказательство. По исходной матрице A строим матрицу

$$B = \frac{1}{\|A\| + \varepsilon} A,$$

где ε — любое положительное число. Тогда

$$\|B\| = \frac{\|A\|}{\|A\| + \varepsilon} < 1$$

и, в силу леммы 2,

$$B^m \xrightarrow{m \rightarrow \infty} O,$$

откуда, согласно лемме 1, следует, что все собственные значения матрицы B меньше единицы по модулю. Но так как собственные значения матрицы B получаются из собственных значений матрицы A умножением на

число $\frac{1}{\|A\|+\varepsilon}$, то для любого собственного значения λ матрицы A должно выполняться неравенство

$$\frac{|\lambda|}{\|A\|+\varepsilon} < 1,$$

или

$$|\lambda| < \|A\| + \varepsilon.$$

Поскольку ε можно взять сколь угодно малым, то $|\lambda| \leq \|A\|$. Лемма доказана.

Теперь уже можно ответить на интересующие нас вопросы относительно сходимости матричной геометрической прогрессии

$$E + A + A^2 + \dots + A^m + \dots$$

Теорема 2. Для того чтобы ряд $E + A + A^2 + \dots + A^m + \dots$ сходиллся, необходимо и достаточно, чтобы $A^m \xrightarrow{m \rightarrow \infty} O$. В этом случае матрица $E - A$ имеет обратную и

$$E + A + A^2 + \dots + A^m + \dots = (E - A)^{-1}.$$

Доказательство. Необходимость этого условия становится очевидной, если вспомнить аналогичное необходимое условие сходимости любого числового ряда и учесть, что сходимость ряда квадратных матриц порядка n эквивалентна сходимости n^2 соответствующих числовых рядов из элементов этих матриц.

Докажем, что условие $A^m \xrightarrow{m \rightarrow \infty} O$ является и достаточным для сходимости ряда $E + A + A^2 + \dots + A^m + \dots$, и найдем сумму этого ряда. Действительно, если $A^m \xrightarrow{m \rightarrow \infty} O$, то по лемме 1 все собственные значения λ_i

($i = 1, 2, \dots, n$) матрицы A по модулю меньше единицы. Тогда все собственные числа матрицы $E - A$, равные $1 - \lambda_i$ ($i = 1, 2, \dots, n$), будут отличны от нуля. Следовательно, определитель этой матрицы, равный произведению всех собственных значений ее, также отличен от нуля, и потому существует матрица $(E - A)^{-1}$.

Рассмотрим тождество

$$(E + A + A^2 + \dots + A^m)(E - A) = E - A^{m+1}.$$

Умножив его справа на матрицу $(E - A)^{-1}$, получим

$$E + A + A^2 + \dots + A^m = (E - A)^{-1} - A^{m+1}(E - A)^{-1}.$$

Так как $A^{m+1} \xrightarrow{m \rightarrow \infty} O$, то

$$E + A + A^2 + \dots + A^m \xrightarrow{m \rightarrow \infty} (E - A)^{-1},$$

т. е.

$$E + A + A^2 + \dots + A^m + \dots = (E - A)^{-1},$$

что и требовалось доказать.

С учетом леммы 1 критерий сходимости бесконечной матричной геометрической прогрессии, данный в теореме 2, может быть сформулирован в другом виде.

Теорема 3. Для того чтобы ряд $E + A + A^2 + \dots + A^m + \dots$ сходил, необходимо и достаточно, чтобы все собственные значения матрицы A были меньше единицы по модулю.

Если же учесть еще и лемму 2 или лемму 3, то можно дать и другой, на этот раз только достаточный признак сходимости рассматриваемой прогрессии.

Теорема 4. Если какая-либо норма матрицы A меньше единицы, то ряд $E + A + A^2 + \dots + A^m + \dots$ сходится.

Последний признак более прост для проверки. При его выполнении нетрудно дать и оценку скорости сходимости рассматриваемого ряда.

Теорема 5. Если $\|A\| < 1$, то

$$\|(E - A)^{-1} - (E + A + A^2 + \dots + A^m)\| \leq \frac{\|A\|^{m+1}}{1 - \|A\|}.$$

Доказательство. Так как при выполнении условия $\|A\| < 1$ ряд $E + A + A^2 + \dots + A^m + \dots$ сходится к матрице $(E - A)^{-1}$, то

$$(E - A)^{-1} - (E + A + A^2 + \dots + A^m) = A^{m+1} + A^{m+2} + \dots$$

и

$$\begin{aligned} \|(E - A)^{-1} - (E + A + A^2 + \dots + A^m)\| &\leq \|A^{m+1}\| + \|A^{m+2}\| + \dots \leq \\ &\leq \|A\|^{m+1} + \|A\|^{m+2} + \dots = \frac{\|A\|^{m+1}}{1 - \|A\|}. \end{aligned}$$

Теорема доказана.

После этих предварительных замечаний приступим, наконец, к рассмотрению методов решения систем линейных алгебраических уравнений.

§ 2.2. ИТЕРАЦИОННЫЕ МЕТОДЫ

Начнем изучение основных методов решения систем линейных алгебраических уравнений с группы итерационных методов. Такие методы, как мы уже отмечали, могут давать точное решение исходной системы, вообще говоря, лишь как результат бесконечного единообразного процесса, называемого процессом итераций. Простота вычислительных схем и однообразие производимых операций делают эти методы удобными при

использовании вычислительной техники. Привлекательным является и свойство самоисправляемости таких методов. Это свойство делает их менее чувствительными по сравнению с точными методами к отдельным ошибкам, допущенным в процессе вычислений. Если при использовании точных методов отдельный сбой в вычислениях неизбежно ведет к ошибке в окончательном результате, то в случае сходящегося итерационного процесса такой сбой влечет за собой, вообще говоря, только лишние приближения. Ошибка, допущенная в каком-то приближении, будет в дальнейшем исправлена последующими приближениями. Однако итерационные методы решения систем линейных алгебраических уравнений не являются универсальными методами. Их сходимость существенным образом зависит от элементов матрицы, определяющей данную задачу. Быстрота сходимости каждого итерационного процесса зависит также и от удачного выбора вектора начального приближения.

2.2.1. Основные разновидности итерационных процессов

Пусть дана система линейных алгебраических уравнений

$$Ax = \bar{f} \quad (2.2.1)$$

с неособенной матрицей A . При построении итерационных методов решения таких систем часто исходную систему (2.2.1) приводят к эквивалентной системе вида

$$\bar{x} = B\bar{x} + \bar{b}. \quad (2.2.2)$$

Тогда последовательность приближений $\bar{x}^{(n)}$ ($n=1, 2, \dots$) к решению $\bar{x}^{(*)}$ этой системы можно строить, например, по рекуррентным формулам

$$\bar{x}^{(k+1)} = B\bar{x}^{(k)} + \bar{b} \quad (k=0, 1, 2, \dots), \quad (2.2.3)$$

при этом начальное приближение $\bar{x}^{(0)}$ можно брать, вообще говоря, произвольным. Систему (2.2.2) можно рассматривать как частный случай операторного уравнения вида $x = \varphi(x)$ и при изучении алгоритмов типа (2.2.3) можно воспользоваться приведенными в гл. 1 результатами исследований метода итерации для таких операторных уравнений. Приведение системы (2.2.1) к виду (2.2.2) можно осуществить по-разному. Например, с помощью любой неособенной матрицы C это преобразование может быть проведено следующим образом:

$$\bar{x} = \bar{x} + C(\bar{f} - A\bar{x}).$$

Здесь $B = E - CA$, $\bar{b} = C\bar{f}$ и алгоритм (2.2.3) принимает вид

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} + C(\bar{f} - A\bar{x}^{(k)}) \quad (k=0, 1, 2, \dots).$$

Если подобные преобразования проводить для каждого шага итераций с новой, вообще говоря, матрицей, то мы придем к алгоритму

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} + C^{(k)}(\bar{f} - A\bar{x}^{(k)}) \quad (k=0, 1, 2, \dots) \quad (2.2.4)$$

или

$$\bar{x}^{(k+1)} = B^{(k)}\bar{x}^{(k)} + \bar{b}^{(k)} \quad (k=0, 1, 2, \dots). \quad (2.2.5)$$

Такой метод итерации обычно называют *нестационарным* в отличие от *стационарного метода* (2.2.3).

Итерационные процессы вида (2.2.4) обладают тем свойством, что точное решение $\bar{x}^{(*)}$ системы (2.2.1) является неподвижной точкой для каждого из них. В самом деле, если в качестве исходного приближения $\bar{x}^{(0)}$ взять вектор $\bar{x}^{(*)}$, то все последующие приближения будут также равны $\bar{x}^{(*)}$. С другой стороны, оказывается, что всякий итерационный процесс приближенного решения системы (2.2.1), для которого $\bar{x}^{(*)}$ является неподвижной точкой, протекающий по формулам (2.2.5), может быть представлен в виде (2.2.4).

Действительно, так как

$$\bar{x}^{(*)} = B^{(k)}\bar{x}^{(*)} + \bar{b}^{(k)},$$

то

$$\begin{aligned} \bar{x}^{(k+1)} &= B^{(k)}\bar{x}^{(k)} + \bar{b}^{(k)} + (\bar{x}^{(*)} - B^{(k)}\bar{x}^{(*)} - \bar{b}^{(k)}) = \bar{x}^{(*)} + B^{(k)}(\bar{x}^{(k)} - \bar{x}^{(*)}) = \\ &= \bar{x}^{(k)} + (B^{(k)} - E)(\bar{x}^{(k)} - \bar{x}^{(*)}) = \bar{x}^{(k)} + (B^{(k)} - E)A^{-1}A(\bar{x}^{(k)} - \bar{x}^{(*)}) = \\ &= \bar{x}^{(k)} + (E - B^{(k)})A^{-1}(\bar{f} - A\bar{x}^{(k)}) = \bar{x}^{(k)} + C^{(k)}(\bar{f} - A\bar{x}^{(k)}), \end{aligned}$$

где

$$C^{(k)} = (E - B^{(k)})A^{-1} \quad (k=0, 1, 2, \dots).$$

При построении итерационных процессов приближенного решения системы (2.2.1) эту систему можно предварительно приводить также и к виду

$$P\bar{x} + Q\bar{x} = \bar{b},$$

где

$$P + Q = CA, \quad \bar{b} = C\bar{f},$$

а выбором неособенной матрицы C мы вправе распорядиться. Аналогично прежнему, и здесь можно построить два типа итерационных алгоритмов: стационарный метод

$$P\bar{x}^{(k+1)} + Q\bar{x}^{(k)} = \bar{b} \quad (k=0, 1, 2, \dots) \quad (2.2.6)$$

и нестационарный метод

$$P^{(k)}\bar{x}^{(k+1)} + Q^{(k)}\bar{x}^{(k)} = \bar{b}^{(k)} \quad (k=0, 1, 2, \dots). \quad (2.2.7)$$

При этом в обоих случаях мы получаем следующее приближение $\bar{x}^{(k+1)}$, вообще говоря, в неявной форме. Поэтому подобные алгоритмы желательно строить так, чтобы матрицу $P^{(k)}$ легко было обратить. Чаще всего берут ее треугольной или диагональной. Соответственно этому методы иногда называют *одношаговыми* в первом случае и *полношаговыми* во втором. В частности, полношаговыми будут все методы типа (2.2.5).

Линейными алгоритмами вида (2.2.5) и (2.2.7) далеко не исчерпываются все разновидности итерационных процессов приближенного решения систем линейных алгебраических уравнений. Вообще говоря, такие процессы могут быть и нелинейными. Например, последовательность приближений к решению системы (2.2.1) можно получать по рекуррентным формулам вида

$$\bar{x}^{(k+1)} = \varphi^{(k)}(\bar{x}^{(0)}, \bar{x}^{(1)}, \dots, \bar{x}^{(k)}) \quad (k=0, 1, 2, \dots), \quad (2.2.8)$$

где $\varphi^{(k)}$ — некоторая функция, зависящая от матрицы системы A , вектора свободных членов \bar{f} , номера приближения k и предыдущих приближений $\bar{x}^{(0)}, \bar{x}^{(1)}, \dots, \bar{x}^{(k)}$.

Мы не станем здесь подробно останавливаться на изучении каждого из типов итерационных процессов, а рассмотрим лишь некоторые из наиболее часто применяемых и характерных итерационных методов.

2.2.2. Метод простой итерации

По-прежнему будем иметь в виду систему (2.2.1). Так как матрица A предполагается неособенной, то решение $\bar{x}^{(*)} = A^{-1}\bar{f}$ этой системы существует и единственно. Будем считать, что исходная система каким-то образом приведена к виду (2.2.2). Пусть также избрано начальное приближение $\bar{x}^{(0)}$ к решению нашей системы. Часто в качестве $\bar{x}^{(0)}$ берут вектор \bar{b} , хотя, вообще говоря, исходное приближение можно выбирать произвольно. Будем называть *методом простой итерации* правило (2.2.3) нахождения последующих приближений к решению нашей системы. Так как правило (2.2.3) линейно, то последовательность приближений $\bar{x}^{(n)}$ ($n=1, 2, \dots$) всегда может быть построена. Если эта последовательность сходится, то она сходится к решению системы (2.2.2). Действительно, если $\bar{x}^{(n)} \xrightarrow{n \rightarrow \infty} \bar{x}^{(*)}$, то предельный переход в алгоритме (2.2.3) приводит нас к равенству $\bar{x}^{(*)} = B\bar{x}^{(*)} + \bar{b}$, что и доказывает сделанное утверждение.

Выясним условия сходимости последовательности приближений, получаемых по методу простой итерации.

Теорема 1. Для того чтобы метод простой итерации (2.2.3) сходиллся при любом начальном приближении $\bar{x}^{(0)}$, необходимо и достаточно, чтобы все собственные значения матрицы B были по модулю меньше единицы.

Доказательство. Проверим сначала достаточность высказанных условий для сходимости метода. Для этого выразим любое приближение, полученное по правилу (2.2.3), через начальное приближение:

$$\begin{aligned}\bar{x}^{(k)} &= B\bar{x}^{(k-1)} + \bar{b} = B(B\bar{x}^{(k-2)} + \bar{b}) + \bar{b} = B^2\bar{x}^{(k-2)} + (E+B)\bar{b} = \\ &= \dots = B^k\bar{x}^{(0)} + (E+B+B^2+\dots+B^{k-1})\bar{b}.\end{aligned}$$

Из этой формулы непосредственно следует проверяемое, ибо

$$B^k \xrightarrow[k \rightarrow \infty]{} O \quad \text{и} \quad E+B+B^2+\dots+B^{k-1} \xrightarrow[k \rightarrow \infty]{} (E-B)^{-1},$$

если все собственные значения матрицы B меньше единицы по модулю.

Пусть теперь при любом $\bar{x}^{(0)}$ существует

$$\lim_{k \rightarrow \infty} \bar{x}^{(k)} = \bar{x}^{(*)}.$$

Тогда

$$\bar{x}^{(*)} = B\bar{x}^{(*)} + \bar{b}.$$

$$\bar{x}^{(*)} - \bar{x}^{(k)} = B(\bar{x}^{(*)} - \bar{x}^{(k-1)}) = B^2(\bar{x}^{(*)} - \bar{x}^{(k-2)}) = \dots = B^k(\bar{x}^{(*)} - \bar{x}^{(0)}),$$

Перейдем в равенстве

$$\bar{x}^{(*)} - \bar{x}^{(k)} = B^k(\bar{x}^{(*)} - \bar{x}^{(0)})$$

к пределу при $k \rightarrow \infty$. Так как вектор $\bar{x}^{(*)} - \bar{x}^{(0)}$ может быть, вообще говоря, любым, а $\bar{x}^{(k)} \xrightarrow[k \rightarrow \infty]{} \bar{x}^{(*)}$, то

$$B^k \xrightarrow[k \rightarrow \infty]{} O,$$

откуда, согласно лемме 1 из § 2.1, следует, что все собственные значения матрицы B меньше единицы по модулю.

Доказанная теорема дает признак сходимости метода простой итерации, который, вообще говоря, трудно проверять, так как связан со спектром матрицы B . Судить о сходимости метода можно и при помощи достаточных признаков, связанных непосредственно с элементами этой

матрицы. Некоторые из таких достаточных признаков вытекают из следующей теоремы.

Теорема 2. Для того чтобы метод простой итерации (2.2.3) сходил, достаточно, чтобы какая-либо норма матрицы B была меньше единицы.

Доказательство. Действительно, если $\|B\| < 1$, то по лемме 3 предыдущего параграфа все собственные значения матрицы B меньше единицы по модулю, и в силу теоремы 1 метод простой итерации (2.2.3) сходится.

На основании последней теоремы можно высказать несколько довольно удобных достаточных признаков сходимости метода простой итерации.

Теорема 3. Метод простой итерации (2.2.3) сходится, если для элементов b_{ij} ($i, j = 1, 2, \dots, n$) матрицы B выполняется одно из следующих условий:

$$1) \sum_{j=1}^n |b_{ij}| < 1 \quad (i=1, 2, \dots, n),$$

$$2) \sum_{i=1}^n |b_{ij}| < 1 \quad (j=1, 2, \dots, n),$$

$$3) \sum_{i,j=1}^n |b_{ij}|^2 < 1.$$

Справедливость сформулированных признаков непосредственно вытекает из теоремы 2, если иметь в виду следующие введенные в п. 2.2.2 нормы матриц:

$$\|B\|_I = \max_{1 \leq i \leq n} \sum_{j=1}^n |b_{ij}|, \quad \|B\|_{II} = \max_{1 \leq j \leq n} \sum_{i=1}^n |b_{ij}|, \quad N(B) = \sqrt{\sum_{i,j=1}^n |b_{ij}|^2}.$$

Можно указать также ряд других достаточных признаков подобного типа, получить которые нетрудно, например, на следующем пути.

Введем в исходной системе линейных алгебраических уравнений

$$x_i = \sum_{j=1}^n b_{ij}x_j + b_i \quad (i=1, 2, \dots, n)$$

новые неизвестные y_i ($i=1, 2, \dots, n$) по формулам

$$x_i = c_i y_i \quad (i=1, 2, \dots, n),$$

где c_i — некоторые (положительные, к примеру) числа.

Тогда исходная система может быть записана в виде

$$y_i = \sum_{j=1}^n b_{ij} \frac{c_j}{c_i} y_j + \frac{b_i}{c_i} \quad (i=1, 2, \dots, n).$$

Будем использовать метод простой итерации для приближенного решения каждой из этих систем, выбрав исходные приближения связанными соотношениями

$$x_i^{(0)} = c_i y_i^{(0)} \quad (i=1, 2, \dots, n).$$

Тогда и компоненты последующих приближений также будут связаны формулами

$$x_i^{(k)} = c_i y_i^{(k)} \quad (i=1, 2, \dots, n; \quad k=1, 2, \dots)$$

и соответствующие процессы простой итерации будут сходиться или расходиться одновременно. Поэтому, опираясь на последнюю теорему, можно утверждать, например, что метод простой итерации (2.2.3) будет сходиться, если найдутся такие положительные числа c_1, c_2, \dots, c_n , что для элементов b_{ij} ($i, j=1, 2, \dots, n$) матрицы B будет выполняться одно из следующих условий:

$$\begin{aligned} 1) \quad & \sum_{j=1}^n |b_{ij}| \frac{c_j}{c_i} < 1 \quad (i=1, 2, \dots, n), \\ 2) \quad & \sum_{i=1}^n |b_{ij}| \frac{c_j}{c_i} < 1 \quad (j=1, 2, \dots, n), \\ 3) \quad & \sum_{i, j=1}^n |b_{ij}|^2 \frac{c_j^2}{c_i^2} < 1. \end{aligned} \tag{2.2.9}$$

Приведенные выше признаки сходимости метода простой итерации позволяют (сравнительно легко в отдельных случаях) получить ответ на вопрос, будет ли сходящимся избранный итерационный процесс. Очень важным для практики является также вопрос о скорости сходимости этого процесса. Получить ответ на этот вопрос помогают *оценки погрешности метода*. Одну из таких оценок дает следующая теорема.

Теорема 4. Если какая-то норма матрицы B , согласованная с данной нормой вектора \bar{x} , меньше единицы, то имеет место следующая оценка погрешности метода простой итерации (2.2.3):

$$\|\bar{x}^{(*)} - \bar{x}^{(k)}\| \leq \|B\|^k \|\bar{x}^{(0)}\| + \frac{\|B\|^k \|\bar{b}\|}{1 - \|B\|}.$$

Доказательство. Как и при доказательстве теоремы 1, находим, что

$$\bar{x}^{(h)} = B^h \bar{x}^{(0)} + (E + B + B^2 + \dots + B^{h-1}) \bar{b}.$$

Так как

$$\|B\| < 1,$$

то

$$\bar{x}^{(*)} = (E + B + B^2 + \dots + B^{h-1} + \dots) \bar{b}$$

и

$$\bar{x}^{(*)} - \bar{x}^{(h)} = (B^h + B^{h+1} + \dots) \bar{b} - B^h \bar{x}^{(0)},$$

откуда

$$\begin{aligned} \|\bar{x}^{(*)} - \bar{x}^{(h)}\| &\leq (\|B\|^h + \|B\|^{h+1} + \dots) \|\bar{b}\| + \|B\|^h \|\bar{x}^{(0)}\| = \\ &= \|B\|^h \|\bar{x}^{(0)}\| + \frac{\|B\|^h \|\bar{b}\|}{1 - \|B\|}, \end{aligned}$$

что и требовалось доказать.

Заметим, что в случае $\bar{x}^{(0)} = \bar{b}$ интересующая нас оценка может быть записана в следующем виде:

$$\|\bar{x}^{(*)} - \bar{x}^{(h)}\| \leq \frac{\|B\|^{h+1} \|\bar{b}\|}{1 - \|B\|}.$$

В самом деле, если за исходное приближение $\bar{x}^{(0)}$ взят вектор-столбец свободных членов \bar{b} , то

$$\bar{x}^{(*)} - \bar{x}^{(h)} = (B^{h+1} + B^{h+2} + \dots) \bar{b}$$

и справедливость сделанного замечания становится очевидной.

Заметим также, что поскольку

$$\bar{x}^{(*)} - \bar{x}^{(h)} = B(\bar{x}^{(*)} - \bar{x}^{(h-1)}),$$

то имеет место и следующее неравенство:

$$\|\bar{x}^{(*)} - \bar{x}^{(h)}\| \leq \|B\| \cdot \|\bar{x}^{(*)} - \bar{x}^{(h-1)}\|.$$

Это неравенство позволяет сравнить точность двух последовательных приближений и часто бывает полезным в практике вычислений.

Полученные выше условия сходимости и оценки погрешности метода простой итерации позволяют теперь уже более целенаправленно подходить и к проблеме выбора такого преобразования исходной системы линейных алгебраических уравнений

$$Ax = \bar{f}$$

к виду

$$\bar{x} = B\bar{x} + \bar{b},$$

удобному для итерации, которое бы обеспечивало сходимость соответствующего процесса простых итераций, при этом с возможно более высокой скоростью. Как мы уже отмечали в предыдущем пункте, такое преобразование может быть проведено, например, с помощью неособенной матрицы C следующим образом:

$$\bar{x} = \bar{x} + C(\bar{f} - A\bar{x}).$$

Здесь роль матрицы B играет матрица $E - CA$, а $\bar{b} = C\bar{f}$. На некоторых наиболее простых или наиболее часто встречающихся способах выбора матрицы C мы сейчас и остановимся.

Прежде всего следует заметить, что матрица $C = A^{-1}$ приводила бы сразу к окончательному решению задачи. Поэтому иногда подбор матрицы C осуществляют путем грубого обращения исходной матрицы A , например, по методу Гаусса (см. § 2.3). Правда, такой подход к выбору матрицы C связан с большим объемом вычислительной работы. Чтобы упростить этот процесс, иногда удобно исходную матрицу A представить предварительно в виде суммы двух таких матриц P и Q , обратная для одной из которых (например, P^{-1}) находится сравнительно просто. Если теперь в качестве матрицы C взять матрицу P^{-1} , то исходная система

$$A\bar{x} = \bar{f}$$

приведется к виду

$$\bar{x} = -P^{-1}Q\bar{x} + P^{-1}\bar{f},$$

удобному для итераций.

В случае, когда матрица A симметрична, можно высказать сравнительно простой критерий сходимости соответствующего итерационного процесса. Прежде чем его сформулировать, заметим, что такое ограничение на матрицу A не является очень обременительным, так как решение системы линейных алгебраических уравнений

$$A\bar{x} = \bar{f}$$

с неособенной матрицей A всегда может быть сведено к решению системы с симметричной и даже положительно определенной матрицей (вещественная симметричная матрица $A = (a_{ij})$ ($i, j = 1, 2, \dots, n$) называется положительно определенной, если квадратичная форма

$$(A\bar{x}, \bar{x}) = \sum_{i,j=1}^n a_{ij}x_i x_j$$

положительно определена, т. е. если все значения этой формы при любых значениях переменных положительны, за исключением значения при $x_1 = x_2 = \dots = x_n = 0$). Такое сведение может быть выполнено с помощью так называемых трансформаций Гаусса, основанных на известной из алгебры теореме о том, что, если A — неособенная матрица, то матрицы $A'A$ и AA' положительно определены.

Оказывается, что, если матрица A симметричная и $A = P + Q$, где P — положительно определенная матрица, обратная для которой известна, то для сходимости метода простой итерации

$$\bar{x}^{(k+1)} = -P^{-1}Q\bar{x}^{(k)} + P^{-1}\bar{f} \quad (k=0, 1, \dots)$$

при любом начальном приближении $\bar{x}^{(0)}$ необходимо и достаточно, чтобы матрицы $P + Q$ и $P - Q$ были положительно определены.

В случае метода Якоби $P=D$, где D — диагональная матрица $[a_{11}, \dots, a_{nn}]$, и сформулированному ранее критерию сходимости здесь можно придать следующую форму.

Для того чтобы метод Якоби для системы $A\bar{x}=\bar{f}$ с симметричной матрицей A , имеющей положительные диагональные элементы, сходиллся при любом выборе начального приближения $x^{(0)}$, необходимо и достаточно, чтобы матрицы A и $2D-A$ (отличающиеся друг от друга знаками недиагональных элементов) были положительно определены.

Условия сходимости метода Якоби можно формулировать, конечно, и через матрицу $B=E-D^{-1}A$, вид которой был выписан выше. Основные результаты здесь даются теоремами 1 и 2. При этом ряду достаточных признаков сходимости метода Якоби, основанных на теореме 2, можно придать несколько более конкретную форму. Так, например, теорема 3 применительно к данному случаю может быть сформулирована в следующем виде.

Метод Якоби для системы $A\bar{x}=\bar{f}$ сходится, если для элементов a_{ij} ($i, j=1, 2, \dots, n$) матрицы A выполняется одно из условий

$$1) \sum'_{j=1}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1 \quad (i=1, 2, \dots, n),$$

$$2) \sum'_{i=1}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1 \quad (j=1, 2, \dots, n),$$

$$3) \sum'_{i,j=1}^n \left| \frac{a_{ij}}{a_{ii}} \right|^2 < 1,$$

где штрихом отмечен тот факт, что при суммировании опускаются слагаемые, отвечающие $i=j$.

Или, если воспользоваться признаками (2.2.9), положив там

$$c_i = \frac{1}{|a_{ii}|} \quad (i=1, 2, \dots, n),$$

то только что выписанные условия (1)–(3) можно заменить соответственно следующими:

$$1) \sum'_{j=1}^n \left| \frac{a_{ij}}{a_{jj}} \right| < 1 \quad (i=1, 2, \dots, n),$$

$$2) \sum'_{i=1}^n \left| \frac{a_{ij}}{a_{jj}} \right| < 1 \quad (j=1, 2, \dots, n);$$

$$3) \sum'_{i,j=1}^n \left| \frac{a_{ij}}{a_{jj}} \right|^2 < 1.$$

Метод Якоби дает быстро сходящийся процесс для случая тех систем линейных алгебраических уравнений, у которых диагональные элементы матрицы системы значительно преобладают по модулю над остальными элементами матрицы. Если же такое доминирование главной диагонали не является значительным, то часто оказывается целесообразным выделять в качестве матрицы P не чисто диагональную матрицу D , а брать, например, матрицу

$$P = \begin{bmatrix} a_{11} & a_{12} & & & 0 \\ a_{21} & a_{22} & & & \\ & & a_{33} & a_{34} & \\ & & a_{43} & a_{44} & \\ 0 & & & & \ddots \end{bmatrix},$$

обращение которой также не представляет особого труда. Роль матрицы C в этом случае играет матрица

$$C = P^{-1} = \begin{bmatrix} \frac{a_{22}}{\Delta_1} & -\frac{a_{12}}{\Delta_1} & & & 0 \\ -\frac{a_{21}}{\Delta_1} & \frac{a_{11}}{\Delta_1} & & & \\ & & 0 & & \\ & & & \ddots & \\ & & & & 0 \end{bmatrix},$$

где обозначено

$$\Delta_1 = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix}.$$

С целью дальнейшего упрощения процесса выбора вспомогательной матрицы C можно выделять из матрицы системы в качестве слагаемого лишь скалярную матрицу

$$P = \alpha E,$$

где α — некоторое отличное от нуля число. Тогда

$$C = \frac{1}{\alpha} E \quad \text{и} \quad B = E - \frac{1}{\alpha} A.$$

При $\alpha = 1$, например, получаем *метод последовательных приближений*

$$\bar{x}^{(k+1)} = B\bar{x}^{(k)} + \bar{b} \quad (k=0, 1, \dots),$$

где

$$B = E - A, \quad \bar{b} = \bar{f}.$$

Выбором константы α иногда можно распорядиться и более рационально. Пусть, например, матрица A положительно определена. В этом случае систему линейных алгебраических уравнений

$$Ax = \bar{f}$$

всегда можно за счет выбора α привести к виду

$$\bar{x} = B\bar{x} + \bar{b}$$

таким образом, чтобы соответствующий процесс простых итераций был сходящимся. В самом деле, поскольку все собственные значения матрицы A положительны и, кроме того (см. лемму 3 из § 2.1), они не превосходят любой из ее норм, то, взяв

$$\alpha = \|A\|,$$

мы построим матрицу

$$B = E - \frac{1}{\alpha} A = E - \frac{1}{\|A\|} A,$$

собственные значения которой будут заведомо удовлетворять условиям теоремы 1 о сходимости метода простой итерации. Если учесть, что решение системы линейных алгебраических уравнений $A\bar{x} = \bar{f}$ с неособенной матрицей A всегда может быть сведено к решению системы с положительно определенной матрицей, то следует отметить не только значительную простоту и эффективность такого подхода к построению алгоритма простых итераций, но и его достаточную универсальность.

Можно говорить и о других способах выбора вспомогательной матрицы C при построении конкретного алгоритма простой итерации. Иногда, например, в качестве матрицы C берут матрицу βA или матрицу $\beta A'$, при этом выбор числового параметра β стараются распорядиться так, чтобы обеспечить сходимость соответствующего алгоритма метода с возможно более высокой скоростью. На этих и других способах построения различных алгоритмов метода простой итерации мы не станем здесь больше останавливаться, а сделаем лишь небольшое замечание, касающееся практической реализации метода.

Практически вычислять простые итерации можно двумя способами. Во-первых, вычисления можно вести непосредственно по формулам

$$\bar{x}^{(k+1)} = B\bar{x}^{(k)} + \bar{b} \quad (k=0, 1, 2, \dots).$$

Здесь каждое найденное приближение можно рассматривать как исходное. Это придает алгоритму простой итерации самоисправляющийся характер. Поэтому на первых шагах процесса нет необходимости, вообще говоря, проводить вычисления с большой точностью: возникающие при этом ошибки впоследствии сглаживаются. Правда, при недостатке опыта такая организация вычислений может быть сопряжена с большим числом лишних итераций.

Во-вторых, k -ю итерацию можно вычислять и по формуле

$$\bar{x}^{(k)} = \bar{b} + B\bar{b} + B^2\bar{b} + \dots + B^k\bar{b},$$

если в качестве исходного приближения $\bar{x}^{(0)}$ взять вектор свободных членов \bar{b} . Здесь вычисления сводятся к нахождению векторов $B\bar{b}$, $B^2\bar{b}$, ..., $B^k\bar{b}$ и последующему их суммированию с вектором \bar{b} . Такая организация вычислений удобна вследствие единообразия процесса, а также потому, что каждое последующее слагаемое является лишь поправкой к найденному приближению. При этом, правда, алгоритм теряет самоисправляющийся характер и становится чувствительным к случайным ошибкам. Кроме того, недостатком этого способа является и возможное накопление

ошибок от округления при возрастании числа слагаемых, что особенно опасно в случае медленно сходящихся процессов.

Чтобы увеличить скорость сходимости метода простой итерации, часто применяют различные приемы ускорения. К рассмотрению некоторых из таких приемов мы еще обратимся в последующих главах, а сейчас отметим лишь, что процесс простой итерации с применением приемов ускорения сходимости в большинстве случаев укладывается в общую схему итерационных методов с нарушением стационарности. К рассмотрению некоторых из нестационарных итерационных процессов мы сейчас и перейдем.

2.2.3. Метод Рундсона

Для приближенного решения исходной системы линейных алгебраических уравнений

$$Ax = \bar{f}$$

будем применять сейчас нестационарные итерационные процессы вида (2.2.5)

$$\bar{x}^{(k+1)} = B^{(k)}\bar{x}^{(k)} + \bar{b}^{(k)} \quad (k=0, 1, 2, \dots).$$

По-прежнему мы будем здесь предполагать, что точное решение $\bar{x}^{(*)} = A^{-1}\bar{f}$ исходной системы (2.2.1) является неподвижной точкой процесса (2.2.5). Тогда, как мы знаем, любой алгоритм типа (2.2.5) может быть записан в виде (2.2.4)

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} + C^{(k)}(\bar{f} - A\bar{x}^{(k)}) \quad (k=0, 1, 2, \dots).$$

При различном выборе матриц $C^{(k)}$ в алгоритмах типа (2.2.4) мы получаем различные итерационные методы. Прежде чем остановиться в качестве примера на способе Рундсона выбора таких матриц, мы выясним общие условия сходимости подобных итерационных процессов.

Так как точное решение системы (2.2.1) является неподвижной точкой алгоритма (2.2.5), то

$$\bar{b}^{(k)} = (E - B^{(k)})A^{-1}\bar{f}$$

и рассматриваемый алгоритм может быть переписан в виде

$$\bar{x}^{(k+1)} = B^{(k)}(\bar{x}^{(k)} - A^{-1}\bar{f}) + A^{-1}\bar{f} \quad (k=0, 1, 2, \dots),$$

откуда непосредственно следует, что

$$\bar{x}^{(k+1)} - A^{-1}\bar{f} = B^{(k)}(\bar{x}^{(k)} - A^{-1}\bar{f}) \quad (2.2.10)$$

для всех $k=0, 1, 2, \dots$. Применяя эту формулу последовательно при $k=0, 1, 2, \dots, n$, получим:

$$\bar{x}^{(n+1)} - A^{-1}\bar{f} = B^{(n)} \cdot B^{(n-1)} \dots B^{(1)} \cdot B^{(0)} (\bar{x}^{(0)} - A^{-1}\bar{f}).$$

Следовательно,

$$\|\bar{x}^{(n+1)} - \bar{x}^{(*)}\| \leq \|B^{(n)}\| \cdot \|B^{(n-1)}\| \dots \|B^{(1)}\| \cdot \|B^{(0)}\| \cdot \|\bar{x}^{(0)} - \bar{x}^{(*)}\|.$$

Если произведение

$$\prod_{k=0}^n \|B^{(k)}\| \xrightarrow{n \rightarrow \infty} 0,$$

то $\|\bar{x}^{(n+1)} - \bar{x}^{(*)}\|$ будет стремиться к нулю при любом выборе исходного приближения $\bar{x}^{(0)}$, т. е.

$$\bar{x}^{(n+1)} \xrightarrow{n \rightarrow \infty} \bar{x}^{(*)}.$$

Для того чтобы рассматриваемое здесь произведение норм матриц стремилось к нулю, достаточно потребовать, чтобы

$$\|B^{(k)}\| \leq \beta < 1$$

для всех $k=0, 1, 2, \dots$. В частности, из этих условий следует известное уже нам достаточное условие сходимости стационарного ($B^{(k)}=B$, $k=0, 1, 2, \dots$) итерационного процесса (2.2.3), даваемое теоремой 2.

Укажем теперь на один из способов выбора матриц $C^{(k)}$ ($k=0, 1, 2, \dots$), при котором соответствующий нестационарный процесс типа (2.2.4) будет заведомо сходящимся.

Будем, например, при построении конкретных итерационных методов вида (2.2.4) выбирать матрицы $C^{(k)}$ скалярными:

$$C^{(k)} = \beta_k E \quad (k=0, 1, 2, \dots).$$

Числовую последовательность β_k ($k=0, 1, 2, \dots$) при этом нужно подобрать так, чтобы сходилась соответствующий итерационный процесс

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} + \beta_k (\bar{f} - A\bar{x}^{(k)}) \quad (k=0, 1, 2, \dots). \quad (2.2.11)$$

Укажем на один из способов построения такой числовой последовательности, например, для случая, когда матрица A положительно определена. В этом случае имеется n положительных собственных значений $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ матрицы A и n соответствующих взаимно ортогональных собственных векторов $\bar{u}^{(1)}, \bar{u}^{(2)}, \dots, \bar{u}^{(n)}$.

Приведем отрезок $[a, b]$ заменой переменной

$$t = \frac{2\lambda - b - a}{b - a}$$

к каноническому отрезку $[-1, 1]$. Тогда многочлен $P_{k+1}(\lambda)$ перейдет в новый многочлен $Q_{k+1}(t)$, при этом

$$Q_{k+1}\left(\frac{a+b}{a-b}\right) = 1,$$

так как

$$P_{k+1}(0) = 1.$$

Известным успехом при решении задачи о минимизации величины M_k было бы построение такого многочлена $Q_{k+1}(t)$ степени $k+1$, принимающего при $t = \frac{a+b}{a-b}$ значение 1, который обладал бы наименьшим максимумом модуля на отрезке $[-1, 1]$. Такой многочлен нетрудно построить с помощью известных многочленов

$$\bar{T}_{k+1}(t) = \frac{1}{2^k} T_{k+1}(t) = \frac{1}{2^k} \cos[(k+1) \arccos t], \quad |t| \leq 1,$$

наименее отклоняющихся от нуля на отрезке $[-1, 1]$.

Дополнительное условие

$$Q_{k+1}\left(\frac{a+b}{a-b}\right) = 1$$

легко учесть, если взять

$$Q_{k+1}(t) = \frac{T_{k+1}(t)}{T_{k+1}\left(\frac{a+b}{a-b}\right)}.$$

Корни многочлена $Q_{k+1}(t)$ совпадают с корнями многочлена Чебышева $T_{k+1}(t)$, которые, как легко видеть, расположены в точках

$$t = \cos \frac{(2i+1)\pi}{2(k+1)} \quad (i=0, 1, 2, \dots, k).$$

Корни же многочлена $P_{k+1}(\lambda)$ расположены в точках

$$\lambda = \frac{1}{\beta_i} \quad (i=0, 1, 2, \dots, k).$$

Учитывая связь между переменными t и λ , можно теперь уже найти и искомые значения β_i :

$$\beta_i = \frac{2}{a+b+(b-a)\cos\frac{(2i+1)\pi}{2(k+1)}} \quad (i=0, 1, \dots, k). \quad (2.2.14)$$

При этом

$$M_k \leq \max_{-1 \leq t \leq 1} |Q_{k+1}(t)| = \frac{1}{\left| T_{k+1}\left(\frac{a+b}{a-b}\right) \right|} < 1, \quad (2.2.15)$$

так как

$$\frac{a+b}{a-b} < -1.$$

Заметим, что подбор чисел β_i ($i=0, 1, \dots, k$) по формуле (2.2.14) можно осуществить лишь в том случае, если их число известно, т. е. если зафиксировано число k . Если же заранее не ясно, сколько шагов итераций потребуется сделать для достижения требуемой точности, то можно использовать числа β_i ($i=0, 1, \dots, k$) циклически, предварительно зафиксировав какое-либо k . При $k=0$ такой процесс будет стационарным, при $k>0$ — нестационарным. Он будет сходящимся в силу неравенств (2.2.13), (2.2.15).

Описанный выше итерационный процесс обычно называют методом Ричардсона.

2.2.4. Метод Зейделя и метод релаксации

В двух предыдущих пунктах настоящего параграфа мы ознакомились с характерными представителями группы полношаговых итерационных методов приближенного решения систем линейных алгебраических уравнений. Сейчас мы остановимся на рассмотрении одношаговых итерационных методов решения таких систем. При этом более детально будет изучен лишь метод Зейделя, являющийся типичным примером стационарного одношагового процесса вида (2.2.6). При рассмотрении же нестационарных методов типа (2.2.7) мы ограничимся только указанием на идею метода релаксации.

При получении алгоритма метода Зейделя мы не станем придерживаться формальной схемы построения стационарных одношаговых процессов, изложенной в п. 2.2.1, а придем к этому методу от известного уже нам метода простой итерации.

Будем считать, что исходная система линейных алгебраических уравнений (2.2.1) каким-то образом уже приведена к виду

$$\bar{x} = B\bar{x} + \bar{b},$$

или

$$x_i = \sum_{j=1}^n b_{ij}x_j + b_i \quad (i=1, 2, \dots, n).$$

Если бы для приближенного решения был избран метод простой итерации, то вычисления мы должны были бы проводить по правилу

$$\bar{x}^{(k+1)} = B\bar{x}^{(k)} + \bar{b} \quad (k=0, 1, 2, \dots)$$

или

$$x_i^{(k+1)} = \sum_{j=1}^n b_{ij}x_j^{(k)} + b_i \quad (i=1, 2, \dots, n; k=0, 1, 2, \dots).$$

При этом алгоритм позволял бы вычислять координаты вектора $\bar{x}^{(k+1)}$ в любом порядке и независимо. Правда, этим самым он лишил бы нас возможности использовать при нахождении последующих координат вектора $\bar{x}^{(k+1)}$ уже найденные координаты этого вектора, хотя последние являются, вообще говоря, улучшенными приближениями к одноименным координатам точного решения $\bar{x}^{(*)}$ по сравнению с соответствующими координатами вектора $\bar{x}^{(k)}$, которые участвуют при вычислениях. Нетрудно изменить алгоритм простой итерации так, чтобы он позволял сразу же использовать при вычислении последующих координат вектора $\bar{x}^{(k+1)}$ уже найденные координаты этого вектора. Например, вычисления можно производить по правилу

$$x_i^{(k+1)} = \sum_{j=1}^{i-1} b_{ij}x_j^{(k+1)} + \sum_{j=i}^n b_{ij}x_j^{(k)} + b_i \quad (2.2.16)$$

$$(i=1, 2, \dots, n; k=0, 1, 2, \dots).$$

Такой итерационный процесс приближенного решения системы линейных алгебраических уравнений и называют методом Зейделя.

Метод Зейделя можно двояко трактовать как разновидность общего итерационного процесса. При первом истолковании за один шаг процесса можно принять переход от вектора

$$(x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})'$$

к вектору

$$(x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_n^{(k+1)})'.$$

В этом случае процесс будет стационарным. При втором истолковании за один шаг процесса можно считать лишь переход от вектора

$$(x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_i^{(k)}, \dots, x_n^{(k)})'$$

к вектору

$$(x_1^{(k+1)}, \dots, x_i^{(k+1)}, x_{i+1}^{(k)}, \dots, x_n^{(k)})'.$$

Такой процесс уже не будет стационарным, но будет циклическим (при первом истолковании за один шаг процесса мы принимали результат применения полного цикла).

В дальнейшем мы будем иметь в виду, как правило, лишь первую трактовку метода.

Если матрицу B разбить на два слагаемых H и F , где

$$H = \begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ b_{21} & 0 & \dots & 0 & 0 \\ b_{31} & b_{32} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nn-1} & 0 \end{bmatrix}, \quad F = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n-1} & b_{1n} \\ 0 & b_{22} & \dots & b_{2n-1} & b_{2n} \\ 0 & 0 & \dots & b_{3n-1} & b_{3n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & b_{nn} \end{bmatrix},$$

то алгоритм (2.2.16) можно переписать в виде

$$\bar{x}^{(k+1)} = H\bar{x}^{(k+1)} + F\bar{x}^{(k)} + \bar{b} \quad (k=0, 1, 2, \dots) \quad (2.2.17)$$

или

$$(E-H)\bar{x}^{(k+1)} = F\bar{x}^{(k)} + \bar{b} \quad (k=0, 1, 2, \dots).$$

Следовательно, метод Зейделя представляет собой одношаговый итерационный процесс вида (2.2.6), где роль треугольной матрицы P играет матрица

$$E-H = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ -b_{21} & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ -b_{n1} & -b_{n2} & \dots & -b_{nn-1} & 1 \end{bmatrix}.$$

Матрица $E-H$ — неособенная, и алгоритм (2.2.17) можно привести к виду

$$\bar{x}^{(k+1)} = (E-H)^{-1}F\bar{x}^{(k)} + (E-H)^{-1}\bar{b} \quad (k=0, 1, 2, \dots).$$

Таким образом, метод Зейделя оказывается эквивалентным методу простой итерации, примененному к системе

$$\bar{x} = (E-H)^{-1}F\bar{x} + (E-H)^{-1}\bar{b},$$

которая равносильна исходной системе линейных алгебраических уравнений. Заметим, что при фактическом проведении процесса вычислительная схема метода Зейделя не совпадает с вычислительной схемой эквивалентного метода простой итерации. Но установленная связь между

этими методами оказывается полезной при выяснении условий сходимости метода Зейделя.

В самом деле, опираясь, например, на теорему 1 о сходимости метода простой итерации, теперь уже можно утверждать, что, для того чтобы метод Зейделя (2.2.17) сходиллся при любом начальном приближении $\bar{x}^{(0)}$, необходимо и достаточно, чтобы все собственные значения матрицы $(E-H)^{-1}F$ были по модулю меньше единицы. Поэтому при выяснении условий сходимости метода Зейделя нас должны интересовать корни уравнения

$$|(E-H)^{-1}F - \lambda E| = 0.$$

Поскольку при построении алгоритма метода Зейделя, как мы видели, в действительности нет необходимости в нахождении матрицы $(E-H)^{-1}F$, то уже даже составление такого уравнения вызывает значительные затруднения. Нетрудно, правда, указать уравнение, корни которого будут совпадать с корнями только что выписанного уравнения, но строиться которое будет более просто. В самом деле, так как определитель произведения квадратных матриц равен произведению определителей этих матриц, а определитель матрицы $E-H$ равен единице, то

$$\begin{aligned} |(E-H)^{-1}F - \lambda E| &= |(E-H)^{-1}(E-H) [(E-H)^{-1}F - \lambda E]| = \\ &= |(E-H)^{-1}| \cdot |F - (E-H)\lambda E| = |F + \lambda H - \lambda E|. \end{aligned}$$

Таким образом, можно высказать следующее утверждение о сходимости метода Зейделя.

Теорема 5. Для того чтобы метод Зейделя (2.2.17) сходиллся при любом начальном приближении $\bar{x}^{(0)}$, необходимо и достаточно, чтобы все корни уравнения

$$|F + \lambda H - \lambda E| = 0$$

были по модулю меньше единицы.

Итак, если исходная система уравнений $A\bar{x} = \bar{f}$ приведена к виду $\bar{x} = B\bar{x} + \bar{b}$, то сходимость метода Зейделя (2.2.17) связана с корнями уравнения

$$\begin{vmatrix} b_{11} - \lambda & b_{12} & \dots & b_{1n} \\ \lambda b_{21} & b_{22} - \lambda & \dots & b_{2n} \\ \dots & \dots & \dots & \dots \\ \lambda b_{n1} & \lambda b_{n2} & \dots & b_{nn} - \lambda \end{vmatrix} = 0.$$

Если же для приближенного решения той же системы уравнений избрать метод простой итерации (2.2.3), то, как мы видели, аналогичную роль играет уравнение

$$\begin{vmatrix} b_{11}-\beta & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22}-\beta & \dots & b_{2n} \\ \cdot & \cdot & \cdot & \cdot \\ b_{n1} & b_{n2} & \dots & b_{nn}-\beta \end{vmatrix} = 0.$$

Уже непосредственное сравнение этих двух уравнений наводит нас на мысль о том, что области сходимости метода простой итерации и метода Зейделя, вообще говоря, различны. Можно привести примеры таких матриц B , для которых метод простой итерации сходится, а метод Зейделя не сходится и наоборот. В самом деле, для случая матрицы

$$B = \begin{bmatrix} 2,5 & -3 \\ 2 & -2,5 \end{bmatrix}$$

уравнение

$$|B - \beta E| = \begin{vmatrix} 2,5-\beta & -3 \\ 2 & -2,5-\beta \end{vmatrix} = \beta^2 - 0,25 = 0$$

имеет корни $\beta_1 = -0,5$ и $\beta_2 = 0,5$, и метод простой итерации будет сходящимся. Метод же Зейделя в случае такой матрицы B сходиться не будет, так как уравнение

$$|F + \lambda H - \lambda E| = \begin{vmatrix} 2,5-\lambda & -3 \\ 2\lambda & -2,5-\lambda \end{vmatrix} = \lambda^2 + 6\lambda - 6,25 = 0$$

имеет один корень, больший единицы по модулю.

Наоборот, для случая матрицы

$$B = \begin{bmatrix} 4,2 & -2 \\ 2 & -0,1 \end{bmatrix}$$

сходиться будет метод Зейделя ($\lambda_1 = -0,6$ и $\lambda_2 = 0,7$), а процесс простой итерации будет расходящимся ($\beta^2 - 4,1\beta + 3,58 = 0$ и $\beta_1\beta_2 > 1$).

Используя установленную связь между методом Зейделя и методом простой итерации и опираясь на теорему 2, можно высказать также и утверждение, дающее одно из достаточных условий сходимости метода Зейделя. А именно, для того чтобы метод Зейделя (2.2.17) сходил, достаточно, чтобы какая-либо норма матрицы $(E - H)^{-1}F$ была меньше единицы. Правда, по упомянутым уже ранее причинам проверка этого условия также затруднительна.

Получим сейчас более просто проверяемые достаточные условия сходимости рассматриваемого метода, которые будут формулироваться непосредственно через элементы матрицы B . Для доказательства справедливости таких условий нам понадобится вспомогательная лемма об определителе матрицы с доминирующими диагональными элементами, которую мы предварительно и докажем.

Лемма 1. Если диагональные элементы матрицы $A = (a_{ij})$ ($i, j = 1, 2, \dots, n$) доминируют по строкам или по столбцам матрицы, т. е. если

$$\sum_{j=1, j \neq i}^n |a_{ij}| < |a_{ii}| \quad (i=1, 2, \dots, n)$$

или

$$\sum_{i=1, i \neq j}^n |a_{ij}| < |a_{jj}| \quad (j=1, 2, \dots, n),$$

то определитель матрицы A отличен от нуля.

Доказательство. Проверим справедливость леммы лишь в случае доминирования диагональных элементов матрицы по строкам (случай доминирования по столбцам исследуется аналогично).

Для доказательства проверяемого утверждения леммы достаточно показать, что система линейных однородных алгебраических уравнений

$$A\bar{x} = \bar{0},$$

где $\bar{x} = (x_1, x_2, \dots, x_n)'$, имеет только тривиальное решение. Предположим противное, т. е. допустим, что эта система имеет и ненулевое решение $\bar{x}^{(*)} = (x_1^{(*)}, x_2^{(*)}, \dots, x_n^{(*)})'$. Среди координат вектора $\bar{x}^{(*)}$ выберем максимальную по модулю:

$$|x_i^{(*)}| \geq |x_j^{(*)}| \quad (j=1, 2, \dots, n).$$

Положим $\bar{x} = \bar{x}^{(*)}$ и рассмотрим соответствующее значение левой части i -го уравнения введенной однородной системы.

Тогда

$$\begin{aligned} |a_{i1}x_1^{(*)} + a_{i2}x_2^{(*)} + \dots + a_{ii}x_i^{(*)} + \dots + a_{in}x_n^{(*)}| &\geq |a_{ii}| \cdot |x_i^{(*)}| - \\ - \sum_{j=1, j \neq i}^n |a_{ij}| \cdot |x_j^{(*)}| &\geq |x_i^{(*)}| \cdot \left(|a_{ii}| - \sum_{j=1, j \neq i}^n |a_{ij}| \right) > 0, \end{aligned}$$

так как по сделанному предположению $|x_i^{(*)}| > 0$, а

$$\sum_{j=1, j \neq i}^n |a_{ij}| < |a_{ii}|$$

по условию леммы. Полученное противоречие и доказывает справедливость высказанного утверждения.

Теперь уже нетрудно доказать и следующую теорему о достаточных условиях сходимости метода Зейделя.

Теорема 6. Для того чтобы метод Зейделя (2.2.16) сходил, достаточно, чтобы выполнялось одно из условий

$$\|B\|_I = \max_{1 \leq i \leq n} \sum_{j=1}^n |b_{ij}| < 1$$

или

$$\|B\|_{II} = \max_{1 \leq j \leq n} \sum_{i=1}^n |b_{ij}| < 1.$$

Доказательство. Рассмотрим только случай первого условия (достаточность условия $\|B\|_{II} < 1$ проверяется аналогично).

Для доказательства выказанного утверждения достаточно показать, что при выполнении условий

$$\sum_{j=1}^n |b_{ij}| < 1 \quad (i=1, 2, \dots, n)$$

значение $\lambda = \lambda^*$, для которого $|\lambda^*| \geq 1$, не может быть корнем уравнения $|F + \lambda H - \lambda E| = 0$ (см. теорему 5). В самом деле, если рассмотреть при таком λ^* сумму модулей недиагональных элементов любой строки определителя

$$|F + \lambda H - \lambda E|,$$

то можно записать:

$$\begin{aligned} & |\lambda^*| \cdot |b_{i1}| + \dots + |\lambda^*| \cdot |b_{ii-1}| + |b_{ii+1}| + \dots + |b_{in}| \leq \\ & \leq |\lambda^*| \sum_{j=1, j \neq i}^n |b_{ij}| = |\lambda^*| \left(\sum_{j=1}^n |b_{ij}| - |b_{ii}| \right) < |\lambda^*| (1 - |b_{ii}|) = \\ & = |\lambda^*| - |\lambda^*| |b_{ii}| \leq |\lambda^*| - |b_{ii}| \leq |\lambda^* - b_{ii}| = |b_{ii} - \lambda^*| \\ & \quad (i=1, 2, \dots, n). \end{aligned}$$

Полученные неравенства

$$\begin{aligned} & |\lambda^*| \cdot |b_{i1}| + \dots + |\lambda^*| \cdot |b_{ii-1}| + |b_{ii+1}| + \dots + |b_{in}| < |b_{ii} - \lambda^*| \\ & \quad (i=1, 2, \dots, n) \end{aligned}$$

представляют собой как раз условия доминирования по строке диагональных элементов матрицы

$$F + \lambda^* H - \lambda^* E.$$

Тогда, согласно лемме 1,

$$|F + \lambda^* H - \lambda^* E| \neq 0,$$

следовательно, при выполнении условия $\|B\|_I < 1$ все корни уравнения

$$|F + \lambda H - \lambda E| = 0$$

по модулю меньше единицы, и метод Зейделя сходится.

Итак, при выяснении вопроса о сходимости данного алгоритма метода Зейделя мы имеем право на основании последней теоремы воспользоваться частью из известных признаков сходимости метода простой итерации.

При практическом использовании метода не менее важно не только выяснить, что данный процесс будет сходящимся, но и знать, как быстро он будет сходиться.

Для метода простой итерации мы имели ряд оценок погрешности, которые позволяли составить представление о скорости сходимости рассматриваемого итерационного процесса. Используя установленную связь между методом простой итерации и методом Зейделя, можно эти оценки перенести и на случай последнего метода. Но использование таких оценок будет затруднено тем, что матрица $(E - H)^{-1}F$ фактически нам не известна. Правда, подобно тому, как это было в только что рассмотренном нами случае с достаточными признаками сходимости, некоторые из оценок погрешности метода простой итерации остаются в силе и для метода Зейделя. Например, мы знаем, что при выполнении условия

$$\|B\|_I = \max_{1 \leq i \leq n} \sum_{j=1}^n |b_{ij}| \leq \mu < 1$$

для метода простой итерации (2.2.3) имеет место оценка

$$\|\bar{x}^{(*)} - \bar{x}^{(k)}\|_I \leq \mu \|\bar{x}^{(*)} - \bar{x}^{(k-1)}\|_I.$$

Оказывается, что в этом случае для метода Зейделя (2.2.16) не только справедлива такая же оценка, но и имеет место несколько лучшая оценка

$$\|\bar{x}^{(*)} - \bar{x}^{(k)}\|_I \leq \mu' \|\bar{x}^{(*)} - \bar{x}^{(k-1)}\|_I,$$

где

$$\mu' = \max_{1 \leq i \leq n} \frac{\gamma_i}{1 - \beta_i}, \quad \beta_i = \sum_{j=1}^{i-1} |b_{ij}|, \quad \gamma_i = \sum_{j=1}^n |b_{ij}|.$$

В самом деле, так как

$$x_i^{(*)} = \sum_{j=1}^n b_{ij} x_j^{(*)} + b_i \quad (i = 1, 2, \dots, n),$$

то, вычитая из этих равенств соответственно равенства

$$x_i^{(k)} = \sum_{j=1}^{i-1} b_{ij} x_j^{(k)} + \sum_{j=i}^n b_{ij} x_j^{(k-1)} + b_i \quad (i=1, 2, \dots, n),$$

получим

$$x_i^{(*)} - x_i^{(k)} = \sum_{j=1}^{i-1} b_{ij} (x_j^{(*)} - x_j^{(k)}) + \sum_{j=i}^n b_{ij} (x_j^{(*)} - x_j^{(k-1)}) \quad (i=1, 2, \dots, n).$$

Тогда

$$\begin{aligned} |x_i^{(*)} - x_i^{(k)}| &\leq \sum_{j=1}^{i-1} |b_{ij}| \cdot |x_j^{(*)} - x_j^{(k)}| + \sum_{j=i}^n |b_{ij}| \cdot |x_j^{(*)} - x_j^{(k-1)}| \leq \\ &\leq \beta_i \|\bar{x}^{(*)} - \bar{x}^{(k)}\|_I + \gamma_i \|\bar{x}^{(*)} - \bar{x}^{(k-1)}\|_I \quad (i=1, 2, \dots, n). \end{aligned}$$

Если $|x_i^{(*)} - x_i^{(k)}|$ достигает максимума при $i=i_0$, то

$$\|\bar{x}^{(*)} - \bar{x}^{(k)}\|_I \leq \frac{\gamma_{i_0}}{1 - \beta_{i_0}} \|\bar{x}^{(*)} - \bar{x}^{(k-1)}\|_I$$

или

$$\|\bar{x}^{(*)} - \bar{x}^{(k)}\|_I \leq \mu' \|\bar{x}^{(*)} - \bar{x}^{(k-1)}\|_I,$$

и утверждаемая оценка доказана.

Осталось только убедиться, что $\mu' \leq \mu$. Действительно, так как при всех $i=1, 2, \dots, n$

$$\beta_i + \gamma_i - \frac{\gamma_i}{1 - \beta_i} = \frac{\beta_i(1 - \beta_i - \gamma_i)}{1 - \beta_i} \geq 0,$$

то

$$\mu' = \max_{1 \leq i \leq n} \frac{\gamma_i}{1 - \beta_i} = \frac{\gamma_{i_1}}{1 - \beta_{i_1}} \leq \beta_{i_1} + \gamma_{i_1} \leq \max_{1 \leq i \leq n} (\beta_i + \gamma_i) \leq \mu,$$

что и требовалось показать.

Однако, как мы видели ранее, метод Зейделя не всегда оказывается более выгодным, чем метод простой итерации. Он даже может расходиться при сходящемся соответствующем процессе простых итераций. Области сходимости этих двух методов, вообще говоря, различны, при этом очень многое здесь зависит от способа приведения исходной системы (2.2.1) к виду (2.2.2), удобному для итерации.

Мы уже знакомились ранее с рядом способов приведения исходной системы линейных алгебраических уравнений к виду (2.2.2). Рассмотрим сейчас лишь один из них, который позволяет для достаточно широкого класса систем построить одношаговый итерационный процесс с более широкой областью сходимости, чем у соответствующего метода простой итерации. Речь пойдет о модификации метода Зейделя, параллельной модификации метода простой итерации, которую мы называли методом Якоби. При этом, как мы знаем, исходная система

$$\sum_{j=1}^n a_{ij}x_j = f_i \quad (i=1, 2, \dots, n)$$

приводится к виду

$$x_i = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j + \frac{f_i}{a_{ii}} \quad (i=1, 2, \dots, n),$$

по которому и записывается алгоритм метода Якоби

$$x_i^{(k+1)} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j^{(k)} - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j^{(k)} + \frac{f_i}{a_{ii}} \quad (i=1, 2, \dots, n; k=0, 1, 2, \dots).$$

Соответствующий одношаговый итерационный процесс

$$x_i^{(k+1)} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j^{(k+1)} - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j^{(k)} + \frac{f_i}{a_{ii}} \quad (i=1, 2, \dots, n; k=0, 1, 2, \dots)$$

часто называют методом Некрасова.

Для этого метода условия сходимости достаточно удобно формулируются посредством исходной матрицы

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}.$$

В самом деле, в случае методов Якоби и Некрасова подготовка системы $A\bar{x} = \bar{f}$ к виду $\bar{x} = B\bar{x} + \bar{b}$ основана, как мы видели, на предварительном умножении системы на диагональную матрицу

$$D^{-1} = [a_{11}, a_{22}, \dots, a_{nn}]^{-1},$$

т. е. здесь

$$B = E - D^{-1}A$$

или

$$B = E - D^{-1}(M + D + N) = -D^{-1}M - D^{-1}N,$$

где

$$M = \begin{bmatrix} 0 & 0 & \dots & 0 \\ a_{21} & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & 0 \end{bmatrix}, \quad N = \begin{bmatrix} 0 & a_{12} & \dots & a_{1n} \\ 0 & 0 & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0 \end{bmatrix}.$$

Таким образом, если придерживаться обозначений, принятых в общем случае метода Зейделя, можно записать, что $H = -D^{-1}M$, а $F = -D^{-1}N$. Матрица $(E - H)^{-1}F$, по собственным значениям которой можно судить о сходимости метода Некрасова, может быть записана в следующем виде:

$$(E - H)^{-1}F = -(E + D^{-1}M)^{-1}D^{-1}N = -[D(E + D^{-1}M)]^{-1}N = -(D + M)^{-1}N.$$

Тогда многочлен $|-(D+M)^{-1}N-\lambda E|$ после умножения на $|-(D+M)|$ принимает вид

$$|N+\lambda(D+M)|.$$

Следовательно, для того чтобы метод Некрасова сходиллся при любом начальном приближении, необходимо и достаточно, чтобы все корни уравнения

$$\begin{vmatrix} a_{11}\lambda & a_{12} & \dots & a_{1n} \\ a_{21}\lambda & a_{22}\lambda & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1}\lambda & a_{n2}\lambda & \dots & a_{nn}\lambda \end{vmatrix} = 0$$

были по модулю меньше единицы.

Если матрица A системы (2.2.1) симметрична, то можно сформулировать еще одно важное условие сходимости метода Некрасова. А именно, для того чтобы метод Некрасова в случае системы линейных алгебраических уравнений с симметричной матрицей A , имеющей положительные диагональные элементы, сходиллся при любом выборе начального приближения, необходимо и достаточно, чтобы матрица A была положительно определена.

Для доказательства достаточности высказанного условия представим симметричную матрицу A в виде

$$A = N' + D + N,$$

где $D = [a_{11}, a_{22}, \dots, a_{nn}]$, а N — треугольная матрица, образованная элементами матрицы A , лежащими выше главной диагонали. Тогда, как мы видели выше, факт сходимости метода Некрасова вполне определяется собственными значениями матрицы $-(D+N')^{-1}N$.

Покажем, что в случае положительно определенной матрицы A все собственные значения матрицы $(D+N')^{-1}N$ по модулю меньше единицы, чем и будет доказана достаточность высказанного условия.

Пусть β — какое-то собственное значение матрицы $(D+N')^{-1}N$, а \bar{z} — соответствующий собственный вектор, т. е.

$$(D+N')^{-1}N\bar{z} = \beta\bar{z}.$$

Тогда

$$N\bar{z} = (D+N')\beta\bar{z} = (A-N)\beta\bar{z} = \beta A\bar{z} - \beta N\bar{z}$$

и

$$(N\bar{z}, \bar{z}) = \beta(A\bar{z}, \bar{z}) - \beta(N\bar{z}, \bar{z}).$$

Отсюда находим

$$\beta = \frac{(N\bar{z}, \bar{z})}{(A\bar{z}, \bar{z}) - (N\bar{z}, \bar{z})}.$$

Обозначим

$$(A\bar{z}, \bar{z}) = p, \quad (D\bar{z}, \bar{z}) = q, \quad (N\bar{z}, \bar{z}) = a + ib.$$

Тогда

$$\beta = \frac{a + ib}{p - a - ib}$$

и

$$|\beta|^2 = \frac{a^2 + b^2}{(p - a)^2 + b^2}.$$

Но

$$p = (A\bar{z}, \bar{z}) = (D\bar{z}, \bar{z}) + (N\bar{z}, \bar{z}) + (N'\bar{z}, \bar{z}) = q + 2a,$$

так как

$$(N'\bar{z}, \bar{z}) = (\bar{z}, N\bar{z}) = a - ib.$$

Поэтому

$$(p-a)^2 = p^2 - 2ap + a^2 = p(p-2a) + a^2 = pq + a^2$$

и

$$|\beta|^2 = \frac{a^2 + b^2}{pq + a^2 + b^2}.$$

Так как $\bar{z} \neq 0$, а матрицы A и D положительно определены, то $p = (A\bar{z}, \bar{z}) > 0$, $q = (D\bar{z}, \bar{z}) > 0$ и $|\beta|^2 < 1$, чем сходимость метода Некрасова доказана.

Докажем теперь необходимость высказанного условия.

Обратим внимание сейчас на циклический характер метода Некрасова и рассмотрим два соседних приближения в $(k+1)$ -м цикле:

$$\bar{\xi}^{(nk+i-1)} = (x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_i^{(k)}, \dots, x_n^{(k)})',$$

$$\bar{\xi}^{(nk+i)} = (x_1^{(k+1)}, \dots, x_i^{(k+1)}, x_{i+1}^{(k)}, \dots, x_n^{(k)})'.$$

Тогда

$$\bar{\xi}^{(nk+i)} = \bar{\xi}^{(nk+i-1)} + E_{ii} D^{-1} (\bar{f} - A \bar{\xi}^{(nk+i-1)}),$$

где

$$E_{ii} = \begin{bmatrix} 0 & \dots & 0 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 1 & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix}^{(i)} (i).$$

Пусть $\bar{x}^{(*)}$ — точное решение системы (2.2.1), а

$$\bar{e}^{(nk+i-1)} = \bar{x}^{(*)} - \bar{\xi}^{(nk+i-1)} \quad \text{и} \quad \bar{e}^{(nk+i)} = \bar{x}^{(*)} - \bar{\xi}^{(nk+i)}$$

соответствующие векторы-ошибки. Очевидно, что

$$\bar{e}^{(nk+i)} = \bar{e}^{(nk+i-1)} - E_{ii} D^{-1} A \bar{e}^{(nk+i-1)} = \bar{e}^{(nk+i-1)} - \frac{1}{a_{ii}} r_i^{(nk+i-1)} \bar{e}_i,$$

где $\bar{e}_i = (0, \dots, 0, 1, 0, \dots, 0)'$, а $r_i^{(nk+i-1)}$ — i -я компонента вектора невязки $\bar{r}^{(nk+i-1)} = \bar{f} - A \bar{\xi}^{(nk+i-1)} = A \bar{e}^{(nk+i-1)}$. Тогда

$$\begin{aligned} (A \bar{e}^{(nk+i)}, \bar{e}^{(nk+i)}) &= (A \bar{e}^{(nk+i-1)}, \bar{e}^{(nk+i-1)}) - 2 \frac{r_i^{(nk+i-1)}}{a_{ii}} (A \bar{e}^{(nk+i-1)}, \bar{e}_i) + \\ &+ \left[\frac{r_i^{(nk+i-1)}}{a_{ii}} \right]^2 (A \bar{e}_i, \bar{e}_i) = (A \bar{e}^{(nk+i-1)}, \bar{e}^{(nk+i-1)}) - \frac{[r_i^{(nk+i-1)}]^2}{a_{ii}}, \end{aligned}$$

так как

$$(A \bar{e}^{(nk+i-1)}, \bar{e}_i) = r_i^{(nk+i-1)}, \quad (A \bar{e}_i, \bar{e}_i) = a_{ii}.$$

Например, не лишена здравого смысла мысль исправлять в первую очередь ту компоненту решения, которая хуже найдена, чтобы при нахождении остальных компонент участвовало уже улучшенное значение ее. Эта идея ослабления влияния «плохой» компоненты может быть осуществлена по-разному. Например, о точности приближенного решения $\bar{x}^{(k)}$ можно судить по величине (в том или ином смысле) вектора ошибки $\bar{\varepsilon}^{(k)} = \bar{x}^{(*)} - \bar{x}^{(k)}$. Правда, этот вектор не может быть вычислен без знания точного решения исходной системы (мы можем лишь оценить его). Иногда вместо вектора $\bar{\varepsilon}^{(k)} = \bar{x}^{(*)} - \bar{x}^{(k)}$ рассматривают вектор $\bar{\delta}^{(k)} = \bar{x}^{(k)} - \bar{x}^{(k-1)}$, который легко вычисляется и по которому в практике вычислений часто судят о близости приближенного решения к точному. Тогда при нахождении вектора $\bar{x}^{(k+1)}$ вычисляют его компоненты в порядке убывания модулей компонент вектора $\bar{\delta}^{(k)} = (\delta_1^{(k)}, \delta_2^{(k)}, \dots, \delta_n^{(k)})'$, а именно, первой находится та компонента вектора $\bar{x}^{(k+1)}$, номер которой совпадает с номером максимальной по модулю компоненты вектора $\bar{\delta}^{(k)}$, потом (с участием только что найденной компоненты) вычисляется та из оставшихся компонент, номер которой совпадает с номером второй по величине модуля среди компонент вектора $\bar{\delta}^{(k)}$, и т. д.

Построенный метод является, очевидно, нестационарным методом вида (2.2.7). Это есть один из примеров так называемых *методов релаксации*.

Принцип релаксации является одним из важных принципов построения итерационных процессов. Он предполагает такой выбор матриц $C^{(k)}$ в алгоритмах типа (2.2.4), например, чтобы на каждом шаге процесса уменьшалась какая-либо величина, характеризующая точность решения исходной системы линейных алгебраических уравнений. Судить о точности решения этой системы можно, скажем, по уже упоминавшемуся вектору ошибки $\bar{\varepsilon}^{(k)} = \bar{x}^{(*)} - \bar{x}^{(k)}$. Подобную же роль может также играть вектор невязки $\bar{r}^{(k)} = \bar{f} - A\bar{x}^{(k)} = A\bar{\varepsilon}^{(k)}$. Методы релаксации могут строиться, например, на уменьшении любой нормы каждого из этих векторов.

Если симметричная матрица A системы (2.2.1) положительно определена, то удобной мерой точности может служить так называемая функция ошибки

$$G(\bar{x}^{(k)}) = (A\bar{\varepsilon}^{(k)}, \bar{\varepsilon}^{(k)}) = (\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) = (A^{-1}\bar{r}^{(k)}, \bar{r}^{(k)}).$$

В силу положительной определенности матрицы A функция ошибки всегда неотрицательна, при этом $G(\bar{x}^{(k)}) = 0$ только при $\bar{x}^{(k)} = \bar{x}^{(*)}$. Так как

$$G(\bar{x}^{(k)}) = (\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) = (\bar{x}^{(*)} - \bar{x}^{(k)}, \bar{f} - A\bar{x}^{(k)}) =$$

$$\begin{aligned}
&= (\bar{x}^{(*)}, \bar{f}) - (\bar{x}^{(k)}, \bar{f}) - (\bar{x}^{(*)}, A\bar{x}^{(k)}) + (\bar{x}^{(k)}, A\bar{x}^{(k)}) = \\
&= (\bar{x}^{(*)}, \bar{f}) - (\bar{x}^{(k)}, \bar{f}) - (\bar{x}^{(k)}, A\bar{x}^{(*)}) + (A\bar{x}^{(k)}, \bar{x}^{(k)}) = \\
&= (A\bar{x}^{(k)}, \bar{x}^{(k)}) - 2(\bar{x}^{(k)}, \bar{f}) + (\bar{x}^{(*)}, \bar{f}),
\end{aligned}$$

то значения функции ошибки лишь постоянным слагаемым отличаются от значений функционала

$$F(\bar{x}^{(k)}) = (A\bar{x}^{(k)}, \bar{x}^{(k)}) - 2(\bar{x}^{(k)}, \bar{f}).$$

Поэтому, хотя функция ошибки и не может быть вычислена без знания точного решения системы $A\bar{x} = \bar{f}$, мы можем судить об убывании функции ошибки, сравнивая между собой соответствующие значения функционала $F(\bar{x}^{(k)})$. Ряд релаксационных методов может быть построен также на принципе уменьшения функции ошибки. Мы не станем здесь приводить примеры таких методов. Несколько методов, основанных на минимизации функционала $F(\bar{x}^{(k)})$ будет построено в § 2.5. Там же будет построен и пример нелинейного метода типа (2.2.8).

§ 2.3. МЕТОДЫ ИСКЛЮЧЕНИЯ

В этом параграфе мы изложим некоторые методы, позволяющие получать точное решение системы линейных алгебраических уравнений в результате выполнения конечного числа арифметических операций. Будут рассмотрены методы, в основе которых лежит идея последовательного исключения неизвестных из уравнений системы. При этом исключение неизвестных из уравнений системы может производиться как путем подходящего комбинирования уравнений системы, так и с помощью специальным образом подобранных матриц (например, матриц вращения, отражения), имеющих целью на одном шаге преобразований обратить в нуль какой-либо элемент матрицы искомой системы или, быть может, обратить в нуль все поддиагональные элементы произвольного столбца этой матрицы. Как в первом, так и во втором случаях преобразования в конечном счете направлены на то, чтобы заданную систему привести к эквивалентной системе и чтобы последняя имела матрицу простого вида.

В сжатой форме большинство подобных методов может быть уложено в следующую схему.

Пусть дана система

$$A\bar{x} = \bar{f}. \quad (2.3.1)$$

Будем преобразовывать эту систему к эквивалентной системе с матрицей простого вида путем умножения ее слева на невырожденные матрицы L_1, L_2, \dots, L_k (о способе выбора таких матриц будет сказано ниже).

$$\left. \begin{aligned} x_m + b_{m\ m+1}x_{m+1} + \dots + b_{mn}x_n &= g_m, \\ a_{m+1\ m+1.m}x_{m+1} + \dots + a_{m+1n.m}x_n &= f_{m+1,m}, \\ \vdots &\quad \vdots \\ a_{n\ m+1.m}x_{m+1} + \dots + a_{nn.m}x_n &= f_{n,m}. \end{aligned} \right\} \quad (2.3.9)$$

где

$$\begin{aligned} b_{mj} &= \frac{a_{mj.m-1}}{a_{mm.m-1}} \quad (j \geq m+1), \quad g_m = \frac{f_{m.m-1}}{a_{mm.m-1}}, \\ a_{ij.m} &= a_{ij.m-1} - a_{im.m-1}b_{mj} \quad (i, j \geq m+1), \\ f_{i.m} &= f_{i.m-1} - a_{im.m-1}g_m. \end{aligned}$$

Предположим, что шаг номера m есть последний возможный шаг преобразований. Могут представиться два случая: $m=n$ и $m<n$. Если $m=n$, то это означает, что после преобразований мы получим систему

$$\left. \begin{array}{l} x_1 + b_{12}x_2 + b_{13}x_3 + \dots + b_{1n}x_n = g_1, \\ x_2 + b_{23}x_3 + \dots + b_{2n}x_n = g_2, \\ \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \\ x_n = g_n \end{array} \right\} \quad (2.3.10)$$

с треугольной матрицей, эквивалентную исходной системе. Из системы (2.3.10) значения для неизвестных находим последовательно от x_n к x_1 по формулам

$$\left. \begin{aligned} x_n &= g_n, \\ x_k &= g_k - b_{kk+1}x_{k+1} - \dots - b_{kn}x_n \quad (k=n-1, n-2, \dots, 1). \end{aligned} \right\} \quad (2.3.11)$$

Процесс нахождения коэффициентов треугольной системы (2.3.10) мы будем называть *прямым ходом*, а процесс получения ее решения по формулам (2.3.11) — *обратным ходом* метода Гаусса.

Пусть $m < n$ и m -е уравнение системы и следующие за ним приведены к виду (2.3.9). Так как шаг m мы считаем последним возможным, то это значит, что в уравнениях (2.3.9), начиная со второго, нельзя выделить ведущего элемента, все $a_{ij,m}$ ($i, j = m+1, \dots, n$) равны нулю и уравнения имеют вид

$$\begin{array}{rcl} x_m + b_{mm+1}x_{m+1} + \dots + b_{mn}x_n & = & g_m, \\ 0 & = & f_{m+1,m}, \\ \cdot & \cdot & \cdot \\ 0 & = & f_{n,m}. \end{array}$$

Если свободные члены $f_{i,m}$ ($i=m+1, \dots, n$) все равны нулю, то получим только одно первое уравнение,

Поэтому данную схему исключения необходимо несколько видоизменить. Достаточно удобной в этом случае будет схема единственного деления с выбором максимального элемента по строке или столбцу или по всей таблице коэффициентов. Например, при выборе максимального элемента по строке в качестве ведущего элемента $(k+1)$ -го шага будем брать максимальный по модулю элемент того уравнения, которое получается из $(k+1)$ -го уравнения после исключения из него всех неизвестных, соответствующих ведущим элементам первых k шагов. Ведущим элементом первого шага будет максимальный по модулю элемент первого уравнения системы (2.3.1).

2.3.2. Метод оптимального исключения

Метод применяется для решения системы (2.3.3) с произвольной неособенной матрицей A . Пусть ведущий элемент первого шага $a_{11} \neq 0$ и это уравнение приведено к виду (2.3.4). Выберем теперь второе уравнение и исключим из него неизвестное x_1 . Мы получим первое уравнение системы (2.3.6):

$$a_{22.1}x_2 + \dots + a_{2n.1}x_n = f_{2.1}.$$

Остальные уравнения системы (2.3.3) оставляем без изменения. Предположим, что $a_{22.1} \neq 0$, и приведем названное уравнение к виду (2.3.8). Исключим из уравнения (2.3.4)

$$x_1 + b_{12}x_2 + \dots + b_{1n}x_n = g_1$$

неизвестное x_2 с помощью уравнения (2.3.8)

$$x_2 + b_{23}x_3 + \dots + b_{2n}x_n = g_{2.1}.$$

Тогда получим

$$\left. \begin{aligned} x_1 + b_{13.1}x_3 + \dots + b_{1n.1}x_n &= g_{1.1}, \\ x_2 + b_{23.1}x_3 + \dots + b_{2n.1}x_n &= g_{2.1}, \end{aligned} \right\} \quad (2.3.13)$$

где

$$\begin{aligned} b_{1j.1} &= b_{1j} - b_{12}b_{2j.1} \quad (j \geq 3), \\ g_{1.1} &= g_1 - b_{12}g_{2.1}, \quad b_{2j.1} = b_{2j}, \quad g_{2.1} = g_2. \end{aligned}$$

Предположим, что после преобразования первых m уравнений система (2.3.3) приведена к эквивалентной системе

$$\left. \begin{aligned} x_1 + b_{1\ m+1\ m-1}x_{m+1} + \dots + b_{1n\ m-1}x_n &= g_{1\ m-1}, \\ x_2 + b_{2\ m+1\ m-1}x_{m+1} + \dots + b_{2n\ m-1}x_n &= g_{2\ m-1}, \\ &\vdots \\ x_m + b_{m\ m+1\ m-1}x_{m+1} + \dots + b_{mn\ m-1}x_n &= g_{m\ m-1}, \\ a_{m+1\ 1}x_1 + \dots + a_{m+1\ m}x_m + a_{m+1\ m+1}x_{m+1} + \dots + a_{m+1\ n}x_n &= f_{m+1}, \\ &\vdots \\ a_{n1}x_1 + \dots + a_{nm}x_m + a_{n\ m+1}x_{m+1} + \dots + a_{nn}x_n &= f_n. \end{aligned} \right\} \quad (2.3.14)$$

Исключим неизвестные x_1, x_2, \dots, x_m из $(m+1)$ -го уравнения этой системы посредством вычитания из него первых k уравнений, умноженных соответственно на $a_{m+1\ 1}, a_{m+1\ 2}, \dots, a_{m+1\ m}$, и разделим вновь полученное уравнение на ведущий элемент $(m+1)$ -го шага (за который мы принимаем коэффициент, стоящий при неизвестном x_{m+1}).

Теперь уравнение примет такой вид:

$$x_{m+1} + b_{m+1\ m+2\ m} x_{m+2} + \dots + b_{m+1\ n\ m} x_n = g_{m+1\ m}.$$

Исключая с помощью этого уравнения неизвестное x_{m+1} из первых m уравнений системы (2.3.14), получим опять систему такого же вида, но с заменой индекса m на $m+1$, при этом:

$$b_{m+1\ p\ m} = \frac{a_{m+1\ p} - \sum_{s=1}^m b_{s\ p\ m-1} a_{m+1\ s}}{a_{m+1\ m+1} - \sum_{s=1}^m b_{s\ m+1\ m} a_{m+1\ s}},$$

$$b_{i\ p\ m} = b_{i\ p\ m-1} - b_{m+1\ p\ m} b_{i\ m+1\ m-1}$$

$$(i=1, 2, \dots, m; \quad p=m+2, m+3, \dots, n),$$

$$g_{m+1\ m} = \frac{f_{m+1} - \sum_{s=1}^m g_{s\ m-1} a_{m+1\ s}}{a_{m+1\ m+1} - \sum_{s=1}^m b_{s\ m+1\ m} a_{m+1\ s}},$$

$$g_{i\ m} = g_{i\ m-1} - g_{m+1\ m} b_{i\ m+1\ m-1}$$

$$(i=1, 2, \dots, m),$$

в предположении, что

$$a_{m+1\ m+1} - \sum_{s=1}^m b_{s\ m+1\ m} a_{m+1\ s} \neq 0.$$

Если все n шагов преобразований возможны, то в результате для искомого решения получим формулы

$$x_i = g_{i\ n-1} \quad (i=1, 2, \dots, n). \quad (2.3.15)$$

Контроль правильности вычислений осуществляется здесь так же, как и в схеме единственного деления. Для решения системы уравнений n -го порядка по методу оптимального исключения необходимо выполнить $\frac{1}{3}n(n^2+3n+2)$ умножений и делений, т. е. почти столько же, сколько и в методе Гаусса.

Метод оптимального исключения по своей структуре весьма близок к методу Гаусса, поэтому его реализация на ЭВМ и реализация метода Гаусса аналогичны. Однако метод оптимального исключения позволяет более эффективно использовать память машины и за счет этого решать системы уравнений приблизительно вдвое большего порядка. Действительно, в силу вида системы (2.3.14) после реализации m -го шага последние $n-m$ уравнений исходной системы остаются без изменения. Поэтому в память машины следует вводить не всю матрицу сразу, а последовательно по одной строке перед каждым шагом. Тогда для проведения $(m+1)$ -го шага достаточно иметь всего

$$\sigma(m) = m(n-m+1) + n + 1$$

рабочих ячеек памяти, которые нужны будут для хранения матрицы

$$\begin{bmatrix} b_{1m+1,m-1} & & b_{1n,m-1} & g_{1,m-1} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ b_{mm+1,m-1} & \dots & b_{mn,m-1} & g_{m,m-1} \end{bmatrix}$$

и коэффициентов $(m+1)$ -го уравнения системы (2.3.14). Так как

$$\max_{1 \leq m \leq n} \sigma(m) = \frac{(n+1)(n+5)}{4} \approx \left(\frac{n}{2}\right)^2,$$

то для решения системы линейных алгебраических уравнений n -го порядка по методу оптимального исключения достаточно иметь поле ячеек величиной $\frac{1}{4}(n+1)(n+5)$, в то время как в методе Гаусса для этой цели необходимо было иметь $(n+1)n$ ячеек. Это позволяет при одинаковом объеме памяти машины решать системы вдвое более высокого порядка, чем по методу Гаусса.

Для осуществления описанной здесь схемы метода оптимального исключения необходимо отличие от нуля всех ведущих элементов. Если этот факт заранее не известен, то целесообразно видоизменить всю схему, перейдя к исключению с выбором главного элемента по строке, как это мы делали в случае схемы единственного деления. Для этого, если в $(m+1)$ -м уравнении после исключения из него x_1, x_2, \dots, x_m максимальным по модулю окажется элемент

$$a_{m+1 p} - \sum_{s=1}^m b_{sp,m-1} a_{m+1 s} \quad (p > m+1),$$

то необходимо переставить местами $(m+1)$ -й и p -й столбцы и продолжить исключение по указанному в методе оптимального исключения правилу.

2.3.3. Метод окаймления

Рассмотрим систему (2.3.1) и предположим матрицу A неособенной. В основе метода окаймления лежит идея вычисления решения системы более высокого порядка через решения вспомогательных систем низшего порядка. Так, например, если

$$B\bar{z}=\bar{b} \quad (2.3.16)$$

есть система линейных алгебраических уравнений некоторого порядка s и

$$G\bar{y}=\bar{g} \quad (2.3.17)$$

есть система линейных алгебраических уравнений порядка $s+1$, в которой матрица G и векторы \bar{y} , \bar{g} определяются по правилу

$$G = \begin{bmatrix} B & \bar{u} \\ \bar{v} & \alpha \end{bmatrix}, \quad \bar{y} = \begin{bmatrix} \bar{\omega} \\ \beta \end{bmatrix}, \quad \bar{g} = \begin{bmatrix} \bar{\tau} \\ \gamma \end{bmatrix},$$

где $\bar{u} = (u_1, u_2, \dots, u_s)'$, $\bar{v} = (v_1, v_2, \dots, v_s)$; α — число; $\bar{\omega}$ и $\bar{\tau}$ — векторы-столбцы размерности s ; β , γ — числа, то между системами указанного вида и их решениями может быть установлена следующая связь:

$$\bar{y} = \begin{bmatrix} \bar{h} \\ 0 \end{bmatrix} - \frac{(\bar{v}, \bar{h}) - \gamma}{(\bar{v}, \bar{f}) - \alpha} \begin{bmatrix} \bar{f} \\ -1 \end{bmatrix}. \quad (2.3.18)$$

В формуле (2.3.18) через \bar{f} обозначено решение системы вида (2.3.16) в случае, когда $\bar{b} = \bar{u}$, а через \bar{h} — аналогичное решение при $\bar{b} = \bar{\tau}$. Действительно, из (2.3.17) получим

$$B\bar{\omega} + \beta\bar{u} = \bar{\tau},$$

$$(\bar{v}, \bar{\omega}) + \alpha\beta = \gamma.$$

Если $\det B \neq 0$, то из первого уравнения находим

$$\bar{\omega} = -\beta B^{-1}\bar{u} + B^{-1}\bar{\tau}. \quad (2.3.19)$$

Умножив полученное выражение скалярно на \bar{v} и учитывая, что

$$(\bar{v}, \bar{\omega}) = \gamma - \alpha\beta,$$

Наряду с системой (2.3.21) рассмотрим систему порядка $k+1$ вида

$$A_{k+1} \bar{x}_{k+1 p} = \bar{b}_{k+1 p}, \quad (2.3.22)$$

где

$$A_{k+1} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1k} & a_{1k+1} \\ a_{21} & a_{22} & \dots & a_{2k} & a_{2k+1} \\ \dots & \dots & \dots & \dots & \dots \\ a_{k1} & a_{k2} & \dots & a_{kk} & a_{kk+1} \\ a_{k+1 1} & a_{k+1 2} & \dots & a_{k+1 k} & a_{k+1 k+1} \end{bmatrix}, \quad \bar{x}_{k+1 p} = (x_{1p}, x_{2p}, \dots, x_{k+1 p})',$$

$$-\bar{b}_{k+1 p} = (a_{1p}, a_{2p}, \dots, a_{k+1 p})', \quad p > k+1.$$

На основании формул (2.3.16)–(2.3.18) установим связь между $\bar{x}_{k+1 p}$ и \bar{x}_{kp} . В нашем случае

$$\bar{u} = -\bar{b}_{kk+1}, \quad \bar{v} = \bar{v}_{k+1} = (a_{k+1 1}, a_{k+1 2}, \dots, a_{k+1 k}), \quad \alpha = a_{k+1 k+1},$$

$$\bar{\tau} = \bar{b}_{kp}, \quad \gamma = -a_{k+1 p}.$$

Следовательно, в силу формулы (2.3.18) для $\bar{x}_{k+1 p}$ получим

$$\bar{x}_{k+1 p} = \begin{bmatrix} \bar{x}_{kp} \\ 0 \end{bmatrix} - \frac{(\bar{v}_{k+1}, \bar{x}_{kp}) + a_{k+1 p}}{(\bar{v}_{k+1}, \bar{x}_{k+1}) - a_{k+1 k+1}} \cdot \begin{bmatrix} \bar{x}_{kk+1} \\ -1 \end{bmatrix}, \quad (2.3.23)$$

$$(1 \leq k \leq n-1, k+1 < p \leq n+1).$$

Заметим здесь, что по идее метода векторы \bar{x}_{kp} и \bar{x}_{kk+1} , являющиеся соответственно решениями систем более низкого порядка, а именно систем

$$A_k \bar{x}_{kp} = \bar{b}_{kp} \quad \text{и} \quad A_k \bar{x}_{kk+1} = \bar{b}_{kk+1},$$

должны быть известны. Таким образом, все вычисления в методе окаймления укладываются в следующую схему:

- 1) сначала вычисляем величины $x_{12}, x_{13}, \dots, x_{1n+1}$ из уравнений

$$a_{11}x_{1p} = b_{1p} \quad (p \geq 2), \quad b_{1p} = -a_{1p}, \quad b_{1n+1} = f_1,$$

в предположении, что $a_{11} \neq 0$;

- 2) затем вычисляем векторы $\bar{x}_{23}, \bar{x}_{24}, \dots, \bar{x}_{2n+1}$ по формуле (2.3.23) при $k=1$ и $3 \leq p \leq n+1$.

Продолжаем этот процесс по аналогии до тех пор, пока не получим вектор \bar{x}_{nn+1} , который, являясь решением системы

$$A_n \bar{x}_{nn+1} = \bar{b}_{nn+1},$$

будет искомым, ибо последняя совпадает с системой (2.3.1).

Рассмотренную схему метода окаймления можно реализовать лишь в том случае, если все коэффициенты

$$(\bar{v}_{k+1}, \bar{x}_{kk+1}) - a_{k+1, k+1} \quad (k=0, 1, \dots, n-1),$$

на которые производится деление в формуле (2.3.18), отличны от нуля. Если окажется, что при некотором значении k коэффициент $(\bar{v}_{k+1}, \bar{x}_{kk+1}) - a_{k+1, k+1}$ равен нулю, то в этом случае целесообразно изменить схему метода окаймления, выполняя в формуле (2.3.23) деление на коэффициент $(\bar{v}_{k+1}, \bar{x}_{ks}) - a_{k+1, s}$, который по абсолютной величине является наибольшим среди всех коэффициентов вида $(\bar{v}_{k+1}, \bar{x}_{kp}) - a_{k+1, p}$ ($k+1 \leq p \leq n$). Тогда формулу (2.3.23) можно переписать так:

$$\bar{x}_{k+1, p} = \begin{bmatrix} \bar{x}_{kp} \\ 0 \end{bmatrix} - \frac{(\bar{v}_{k+1}, \bar{x}_{kp}) + a_{k+1, p}}{(\bar{v}_{k+1}, \bar{x}_{ks}) - a_{k+1, s}} \begin{bmatrix} \bar{x}_{ks} \\ -1 \end{bmatrix} \quad (2.3.24)$$

$$(1 \leq k \leq n-1, p = k+1, k+2, \dots, s-1, s+1, \dots, n+1).$$

Определяемую этой формулой схему метода окаймления называют схемой с выбором максимального элемента по строке.

Укажем в заключение на связь метода окаймления с методом оптимального исключения. Рассмотрим вектор

$$\bar{b}_{p, m-1} = (b_{1p, m-1}, b_{2p, m-1}, \dots, b_{mp, m-1})',$$

компонентами которого являются элементы p -го столбца матрицы системы (2.3.14). Сравнивая правило получения $\bar{b}_{p, m-1}$ и вектора \bar{x}_{mp} по формуле (2.3.23), убедимся в том, что

$$\bar{x}_{mp} = \bar{b}_{p, m-1} \quad (p = m+1, m+2, \dots, n).$$

Как и в методе оптимального исключения для решения системы линейных алгебраических уравнений n -го порядка по методу окаймления достаточно иметь поле рабочих ячеек величиной

$$\max_{1 \leq m \leq n} \sigma(m) = \frac{(n+1)(n+5)}{4},$$

ибо на $(m+1)$ -м шаге при вычислении векторов $\bar{x}_{m+1, p}$ ($p > m+1$) тре-

буется знание только векторов \bar{x}_{mp} ($p > m$) и коэффициентов $(m+1)$ -го уравнения исходной системы, т. е. числового массива величиной

$$\sigma(m) = m(n-m+1) + n + 1.$$

Для решения системы (2.3.1) по методу окаймления необходимо выполнить $\frac{1}{6}n(2n^2+9n+1)$ умножений и делений — примерно столько же, сколько в методах Гаусса и оптимального исключения.

2.3.4. Вычисление определителей

Каждая из рассмотренных в пп. 2.3.1—2.3.3 схем для решения систем может быть применена и для вычисления определителей. Остановимся сначала на описании схемы единственного деления. Пусть

$$\Delta = \begin{vmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{vmatrix}$$

и пусть $a_{11} \neq 0$. Вынося элемент a_{11} из первой строки, получим

$$\Delta = a_{11} \begin{vmatrix} 1 & b_{12} & \dots & b_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{vmatrix},$$

где величины b_{1j} определяются по формулам (2.3.5). Вычитая из каждой строки, начиная со второй, первую строку, умноженную соответственно на $a_{21}, a_{31}, \dots, a_{n1}$, мы получим, очевидно,

$$\Delta = a_{11} \cdot \begin{vmatrix} 1 & b_{12} & \dots & b_{1n} \\ 0 & a_{22.1} & \dots & a_{2n.1} \\ \dots & \dots & \dots & \dots \\ 0 & a_{n2.1} & \dots & a_{nn.1} \end{vmatrix} = a_{11} \begin{vmatrix} a_{22.1} & a_{23.1} & \dots & a_{2n.1} \\ a_{32.1} & a_{33.1} & \dots & a_{3n.1} \\ \dots & \dots & \dots & \dots \\ a_{n2.1} & a_{n3.1} & \dots & a_{nn.1} \end{vmatrix},$$

где величины $a_{ij.1}$ определяются по формулам (2.3.7). С образовавшимся определителем $(n-1)$ -го порядка поступаем совершенно таким же образом, если только $a_{22.1} \neq 0$.

Продолжая процесс, мы получим, что искомый определитель равен произведению ведущих элементов:

$$\Delta = a_{11} \cdot a_{22.1} \cdot a_{33.2} \dots a_{nn.n-1}.$$

Аналогично вычисляется определитель и в методе оптимального исключения, если все ведущие элементы отличны от нуля. При этом имеем:

$$\Delta = a_{11} \prod_{m=1}^{n-1} \left(a_{m+1, m+1} - \sum_{s=1}^m b_{sm+1, m} a_{m+1s} \right),$$

Если хотя бы один ведущий элемент равен нулю и схема метода оптимального исключения реализуется с выбором максимального элемента по строке, то определитель также будет равен произведению этих новых ведущих элементов, которые мы обозначим через α_k ($k=1, 2, \dots, n$). Однако в этом случае, чтобы сохранить знак определителя, надо каждый элемент α_k умножить на $(-1)^{l_k+1}$, где l_k — номер неизвестного, исключенного на $(k+1)$ -м шаге, если все неизвестные, не исключенные на первых k шагах, были занумерованы подряд слева направо числами $1, 2, \dots, n-k$.

Таким образом,

$$\Delta = \prod_{k=1}^n (-1)^{l_k+1} \cdot \alpha_k.$$

При вычислении определителя по любой из этих трех формул может для некоторого числа сомножителей произойти переполнение разрядной сетки машины (или образование машинного нуля), хотя сам определитель Δ не очень велик (мал). Этого можно избежать, вычисляя определитель, например, по третьей формуле так:

$$\Delta = \left(q \prod_j (-1)^{l_j+1} \cdot \alpha_j \right) \left(r \prod_k (-1)^{l_k+1} \cdot \alpha_k \right),$$

где q близко к максимальному допустимому в ЭВМ числу, r — близко к минимальному, причем $qr=1$ и все $|\alpha_j| \leq 1$, а $|\alpha_k| \geq 1$.

Остановимся, наконец, на правиле вычисления определителя в методе окаймления. Для этого установим прежде всего связь между определителями матриц B и

$$G = \begin{bmatrix} B & \bar{u} \\ \bar{v} & \alpha \end{bmatrix}.$$

При описании метода окаймления мы рассматривали систему

$$\begin{aligned} B\bar{\omega} + \beta\bar{u} &= \bar{\tau}, \\ (\bar{v}, \bar{\omega}) + \beta\alpha &= \gamma. \end{aligned}$$

Отсюда, используя выражение для $\bar{\omega}$, можно перейти к эквивалентной системе

$$\begin{aligned} B\bar{\omega} + \beta\bar{u} &= \bar{\tau}, \\ \beta[\alpha - (\bar{v}, B^{-1}\bar{u})] &= \gamma - (\bar{v}, B^{-1}\bar{\tau}), \end{aligned}$$

определитель которой по теореме Лапласа равен

$$(\alpha - (\bar{v}, B^{-1}\bar{u})) \cdot |B|.$$

Следовательно,

$$|G| = (\alpha - (\bar{v}, B^{-1}\bar{u})) |B|. \quad (2.3.25)$$

Теперь по аналогии устанавливается связь между определителями матриц A_k и A_{k+1} , а именно:

$$\begin{aligned} |A_{k+1}| &= (a_{k+1\ k+1} - (\bar{v}_{k+1}, \bar{x}_{k\ k+1})) \cdot |A_k| \\ &\quad (k=1, 2, \dots, n-1). \end{aligned} \quad (2.3.26)$$

Применяя рекуррентно формулу (2.3.26) для определителя матрицы A , получим

$$\Delta = a_{11} \prod_{k=1}^{n-1} (a_{k+1\ k+1} - (\bar{v}_{k+1}, \bar{x}_{k\ k+1})). \quad (2.3.27)$$

Отметим здесь, что из формулы (2.3.26) следует, что множители $(a_{k+1\ k+1} - (\bar{v}_{k+1}, \bar{x}_{k\ k+1}))$ будут все отличными от нуля тогда и только тогда, когда все главные миноры матрицы A будут отличными от нуля. Если этот факт известен заранее ($|A_k| \neq 0$), то метод окаймления для системы $A\bar{x} = \bar{f}$ можно реализовать по схеме, определяемой формулой (2.3.23).

Сравнивая процесс Гаусса для решения системы с процессом вычисления определителя, мы видим, что объем вычислений для решения системы лишь немногим превышает объем вычисления одного определителя. Этим, в частности, объясняется то, что пользоваться формулами Крамера для численного решения системы не целесообразно.

2.3.5. Обращение матриц

Задача решения системы линейных алгебраических уравнений тесно связана с задачей обращения матрицы, поэтому все рассмотренные выше методы исключения можно приспособить также и к нахождению обратной матрицы. Действительно, если матрица A неособенна, то по определению обратной матрицы

$$AX=E, \quad (2.3.28)$$

где через X обозначена обратная к A матрица. Пусть $\bar{x}_k = (x_{k1}, x_{k2}, \dots, x_{kn})'$ — k -й столбец матрицы X и $\bar{e}_k = (0, 0, \dots, 0, 1, 0, \dots, 0)'$ — единич-

ный вектор. Тогда в силу (2.3.28) следует, что определение элементов обратной матрицы эквивалентно решению n систем линейных алгебраических уравнений вида

$$A\bar{x}_k = \bar{e}_k \quad (k=1, 2, \dots, n).$$

Для контроля вычисления и оценки точности результата целесообразно произвести умножение A на $X=A^{-1}$.

С другой стороны, если известна A^{-1} , то сразу можно записать решение \bar{x} любой системы $A\bar{x} = \bar{f}$ в виде

$$\bar{x} = A^{-1}\bar{f}.$$

Метод оптимального исключения может быть более эффективно применен к задаче обращения матрицы. Ниже мы изложим один из вариантов метода исключения для обращения матрицы, принадлежащий Жордану. Суть его в следующем. Пусть требуется найти матрицу, обратную к A , и пусть матрица A приведена к виду

$$A_k = \begin{bmatrix} 1 & 0 & \dots & 0 & a_{1\ k+1}^{(k)} & a_{1n}^{(k)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & & 1 & a_{k\ k+1}^{(k)} & a_{kn}^{(k)} \\ 0 & 0 & & 0 & a_{k+1\ k+1}^{(k)} & a_{k+1n}^{(k)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & & 0 & a_{n\ k+1}^{(k)} & a_{nn}^{(k)} \end{bmatrix}, \quad (2.3.29)$$

Предположим, что $a_{k+1\ k+1}^{(k)} \neq 0$, и разделим $(k+1)$ -ю строку матрицы A_k на этот коэффициент, который назовем ведущим элементом. С помощью этой преобразованной строки исключим все внедиагональные элементы $(k+1)$ -го столбца, для чего будем последовательно умножать эту строку

на $a_{1\ k+1}^{(k)}, \dots, a_{k\ k+1}^{(k)}, a_{k+2\ k+1}^{(k)}, \dots, a_{n\ k+1}^{(k)}$ и вычитать соответственно из первой, второй и т. д. строк. После выполнения этих операций мы придем к матрице A_{k+1} вида

$$A_{k+1} = \begin{bmatrix} 1 & 0 & \dots & 0 & a_{1k+2}^{(k+1)} & \dots & a_{1n}^{(k+1)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & 1 & a_{k+1, k+2}^{(k+1)} & \dots & a_{k+1n}^{(k+1)} \\ 0 & 0 & \dots & 0 & a_{k+2, k+2}^{(k+1)} & \dots & a_{k+2n}^{(k+1)} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & & 0 & a_{nk+2}^{(k+1)} & \dots & a_{nn}^{(k+1)} \end{bmatrix},$$

которая, как это легко проверить, связана с A_k равенством

$$A_{k+1} = L_{k+1} A_k, \quad (2.3.30)$$

где

$$L_{k+1} = \begin{bmatrix} 1 & \dots & 0 & -\frac{a_{1k+1}^{(k)}}{a_{k+1, k+1}^{(k)}} & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & & 1 & -\frac{a_{hk+1}^{(k)}}{a_{k+1, k+1}^{(k)}} & \dots & 0 \\ 0 & & 0 & \frac{1}{a_{k+1, k+1}^{(k)}} & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \underbrace{\dots}_k & 0 & -\frac{a_{nk+1}^{(k)}}{a_{k+1, k+1}^{(k)}} & \dots & 1 \end{bmatrix}.$$

Преобразуя последовательно матрицу $A_0 = A$ матрицами L_1, L_2, \dots, L_n , мы получим матрицы A_1, A_2, \dots, A_n , причем $A_n = E$. Таким образом,

$$\begin{aligned} A_0 &= A, \\ A_1 &= L_1 A_0, \\ A_2 &= L_2 A_1, \\ &\cdot \quad \cdot \quad \cdot \\ A_n &= E = L_n A_{n-1}. \end{aligned}$$

Отсюда

$$E = L_n \cdot L_{n-1} \dots L_1 \cdot A$$

и, значит,

$$A^{-1} = L_n \cdot L_{n-1} \dots L_1. \quad (2.3.31)$$

Матрицу A^{-1} удобно вычислять по следующей рекурсионной формуле:

$$B_{k+1} = L_{k+1} B_k \quad (k=1, 2, \dots, n-1), \quad (2.3.32)$$

где полагаем $B_1 = L_1$. Очевидно, что $B_n = A^{-1}$. В силу формулы (2.3.32) переход от матрицы B_k к матрице B_{k+1} осуществляется таким же образом, как и переход от A_k к A_{k+1} .

Отметим еще, что матрицы A_k и B_k имеют специальный вид: в A_k первые k столбцов, а в B_k последние $n-k$ столбцов совпадают с соответствующими столбцами единичной матрицы и не нуждаются в отдельном запоминании их в памяти машины. Поэтому обе матрицы можно хранить и преобразовывать на том же месте, где хранилась матрица A .

Обозначим $C_0 = A$ и через C_k обозначим квадратную матрицу порядка n , первыми k столбцами которой являются столбцы матрицы B_k (начиная с первого и до k -го включительно). Оставшиеся $n-k$ столбцов матрицы C_k являются соответственно (начиная с $(k+1)$ -го столбца) столбцами матрицы A_k . Тогда вычисление обратной матрицы, учитывая формулы (2.3.30) — (2.3.32), сводится к построению последовательности матриц C_1, C_2, \dots, C_n , коэффициенты которых определяются по таким рекуррентным формулам:

$$c_{ij}^{(k+1)} = \begin{cases} c_{ij}^{(k)} - c_{k+1j}^{(k)} \frac{c_{ik+1}^{(k)}}{c_{k+1k+1}^{(k)}}, & i \neq k+1, \quad j \neq k+1, \\ \frac{c_{k+1j}^{(k)}}{c_{k+1k+1}^{(k)}}, & i = k+1, \quad j \neq k+1, \\ -\frac{c_{ik+1}^{(k)}}{c_{k+1k+1}^{(k)}}, & i \neq k+1, \quad j = k+1, \\ \frac{1}{c_{k+1k+1}^{(k)}}, & i = k+1, \quad j = k+1, \end{cases} \quad (2.3.33)$$

$$(i, j = 1, 2, \dots, n; \quad k = 0, 1, 2, \dots, n-1).$$

Таким образом, при обращении матрицы A по методу Жордана мы используем только формулы (2.3.33) и окончательный результат получаем в виде матрицы C_n , которая равна A^{-1} . При этом необходимо выполнить n^3 умножений и делений.

Если среди главных миноров матрицы A есть равные нулю, то тогда нулю может равняться и какой-либо ведущий элемент. В этом случае целесообразно применить схему метода Жордана с выбором максимального элемента по строке. Тогда в качестве ведущего элемента на $(k+1)$ -м

шаге берется тот коэффициент в $(k+1)$ -й строке, который среди ее коэффициентов, находящихся в столбцах, не исключенных на первых k шагах, является максимальным по модулю. Перед выполнением преобразований (2.3.33) целесообразно переставить $(k+1)$ -й столбец и столбец с ведущим элементом. Чтобы сохранить матрицу A^{-1} , надо в матрице C_n переставить строки в порядке, обратном порядку перестановки столбцов.

Рассмотрим еще применение идеи метода окаймления к задаче нахождения обратной матрицы. Пусть требуется найти обратную к A матрицу и пусть A_k — некоторая неособенная квадратная матрица, для которой обратная матрица известна. Установим связь между матрицами A_k^{-1} и A_{k+1}^{-1} в предположении, что

$$A_{k+1} = \begin{bmatrix} A_k & \bar{u}_k \\ \bar{v}_k & \alpha_k \end{bmatrix}$$

и $\bar{u}_k, \bar{v}_k, \alpha_k$ — известные векторы и число. Пусть

$$A_{k+1}^{-1} = \begin{bmatrix} P_k & \bar{r}_k \\ \bar{q}_k & \frac{1}{\beta_k} \end{bmatrix}.$$

Нам надо, считая матрицу A_k^{-1} известной, найти матрицу A_{k+1}^{-1} , т. е. определить матрицу P_k , вектор-столбец \bar{r}_k , вектор-строку \bar{q}_k и число $\frac{1}{\beta_k}$. По определению обратной матрицы

$$A_{k+1}A_{k+1}^{-1} = \begin{bmatrix} A_k P_k + \bar{u}_k \bar{q}_k & A_k \bar{r}_k + \bar{u}_k \frac{1}{\beta_k} \\ \bar{v}_k P_k + \alpha_k \bar{q}_k & (\bar{v}_k \bar{r}_k) + \alpha_k \frac{1}{\beta_k} \end{bmatrix} = \begin{bmatrix} E_k & 0 \\ 0 & 1 \end{bmatrix},$$

где E_k — единичная матрица порядка k . Отсюда получим:

$$\left. \begin{aligned} A_k P_k + \bar{u}_k \bar{q}_k &= E_k, \\ A_k \bar{r}_k + \frac{1}{\beta_k} \bar{u}_k &= 0, \\ \bar{v}_k P_k + \alpha_k \bar{q}_k &= 0, \\ (\bar{v}_k, \bar{r}_k) + \alpha_k \frac{1}{\beta_k} &= 1. \end{aligned} \right\} \quad (2.3.34)$$

По существу система (2.3.34) является системой $(k+1)^2$ -линейных алгебраических уравнений для определения такого же числа неизвестных — элементов обратной матрицы A_{k+1}^{-1} . Для наших целей систему (2.3.34) удобно рассматривать как систему уравнений с четырьмя неизвестными: P_k , \bar{q}_k , \bar{r}_k , $\frac{1}{\beta_k}$, ибо по условию нам известна матрица A_k^{-1} . Поэтому из второго уравнения системы (2.3.34) имеем

$$\bar{r}_k = -\frac{1}{\beta_k} A_k^{-1} \bar{u}_k. \quad (2.3.35)$$

Подставив это выражение в четвертое уравнение системы (2.3.34), получим

$$\left(\bar{v}_k, -\frac{1}{\beta_k} A_k^{-1} \bar{u}_k \right) + \frac{\alpha_k}{\beta_k} = 1.$$

Следовательно,

$$\beta_k = \alpha_k - (\bar{v}_k, A_k^{-1} \bar{u}_k). \quad (2.3.36)$$

Из первого уравнения системы (2.3.34) определим матрицу P_k :

$$P_k = A_k^{-1} - A_k^{-1} \bar{u}_k \bar{q}_k. \quad (2.3.37)$$

Наконец, из третьего уравнения системы (2.3.34), используя формулу (2.3.37), найдем

$$\begin{aligned} \bar{v}_k (A_k^{-1} - A_k^{-1} \bar{u}_k \bar{q}_k) + \alpha_k \bar{q}_k &= \bar{v}_k A_k^{-1} + \\ + (\alpha_k - (\bar{v}_k, A_k^{-1} \bar{u}_k)) \bar{q}_k &= \bar{v}_k A_k^{-1} + \beta_k \bar{q}_k = \bar{0} \end{aligned}$$

и отсюда получим \bar{q}_k :

$$\bar{q}_k = -\frac{\bar{v}_k A_k^{-1}}{\beta_k}. \quad (2.3.38)$$

Учитывая это выражение для \bar{q}_k , из (2.3.37) для матрицы P_k окончательно получим такое выражение:

$$P_k = A_k^{-1} + \frac{A_k^{-1} \bar{u}_k \bar{v}_k A_k^{-1}}{\beta_k}. \quad (2.3.39)$$

Таким образом, зная обратную матрицу A_k^{-1} , мы сможем с помощью формул (2.3.35) — (2.3.39) вычислить обратную матрицу к A_{k+1} , при этом

$$A_{k+1}^{-1} = \begin{bmatrix} A_k^{-1} + \frac{A_k^{-1} \bar{u}_k \bar{v}_k A_k^{-1}}{\beta_k} & -\frac{A_k^{-1} \bar{u}_k}{\beta_k} \\ -\frac{\bar{v}_k A_k^{-1}}{\beta_k} & \frac{1}{\beta_k} \end{bmatrix}. \quad (2.3.40)$$

Характерно, что при вычислении матрицы A_{k+1}^{-1} по формуле (2.3.40) мы нигде не занимаемся обращением матриц как таковым, а выполняем только такие более простые операции, как умножение матриц, умножение матрицы на вектор и деление на число β_k (при вычислении β_k выполняется операция скалярного умножения векторов). Заметим здесь, что в силу формулы (2.3.25) число β_k имеет следующий смысл:

$$\beta_k = \frac{|A_{k+1}|}{|A_k|},$$

ибо применительно к матрице A_{k+1} эту формулу можно записать в виде

$$|A_{k+1}| = (\alpha_k - (\bar{v}_k, A_k^{-1} \bar{u}_k)) \cdot |A_k|.$$

Как описанный здесь процесс применить к решению задачи об обращении матрицы A ? Предположим, что эта матрица имеет отличные от нуля главные миноры. Тогда, последовательно обращая по формуле (2.3.40) матрицы

$$A_1 = [a_{11}], \quad A_2 = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \quad A_3 = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}, \quad \dots,$$

каждая из которых получается окаймлением предыдущей, найдем матрицу A^{-1} . Для этого необходимо выполнить n^3 умножений и делений.

Особенно метод окаймления эффективен при обращении эрмитовых и треугольных матриц. Действительно, если A эрмитова, то для всех k в этом случае имеют место равенства

$$A_{k+1}^{-1} = (A_{k+1}^{-1})^* \quad \text{и} \quad \bar{v}_k = \bar{u}_k^*,$$

ибо

$$(A_{k+1}^* A_{k+1}^{-1})^* = E = (A_{k+1}^{-1})^* A_{k+1}.$$

Следовательно, матрицу A_{k+1}^{-1} можно записать в таком виде:

$$A_{k+1}^{-1} = \begin{bmatrix} A_k^{-1} + \frac{\bar{p}_k \bar{p}_k^*}{\beta_k} & -\frac{\bar{p}_k}{\beta_k} \\ -\frac{\bar{p}_k^*}{\beta_k} & \frac{1}{\beta_k} \end{bmatrix},$$

где

$$\bar{p}_k = A_k^{-1} \bar{u}_k, \quad \beta_k = \alpha_k - (\bar{u}_k^*, \bar{p}_k).$$

Поскольку в эрмитовых матрицах после определения вектора $\bar{q}_k = -\frac{\bar{p}_k}{\beta_k}$ вектор \bar{r}_k определяется путем перехода к транспонированному и комплексно сопряженному вектору \bar{q}_k , то для обращения эрмитовой матрицы необходимо выполнить вдвое меньше арифметических операций, чем для обращения произвольной матрицы, например надо выполнить $\frac{1}{2}n^2(n+1)$ умножений и делений.

Еще более простым будет алгоритм обращения треугольной матрицы. Рассмотрим, например, правую треугольную матрицу A . Для нее все $\bar{v}_k = 0$, а A_k^{-1} — треугольные матрицы того же наименования. Поэтому в силу формулы (2.3.40)

$$A_{k+1}^{-1} = \begin{bmatrix} A_k^{-1} & -\frac{A_k^{-1} \bar{u}_k}{\beta_k} \\ 0 & \frac{1}{\beta_k} \end{bmatrix}.$$

Для обращения треугольной матрицы необходимо выполнить $\frac{1}{6}n(n^2 + 3n + 2)$ умножений и делений.

Понятно, что описанные схемы метода окаймления проходят лишь для матриц с отличными от нуля главными минорами. В общем случае надо применять схему с выбором главного элемента по строке аналогично тому, как мы это делали в методе Жордана.

Обращение матрицы по указанным выше схемам не дает уверенности в точности полученных результатов из-за неизбежных округлений, влияние которых на конечный результат трудно оценить. Поэтому, как мы уже отмечали, для контроля точности вычисления обратной матрицы надо выполнить умножение матрицы на ее обратную и результат сравнить с единичной матрицей. Несовпадение произведения с единичной

матрицей указывает на степень неточности в вычислении обратной матрицы. Если неточности большие и результат нельзя считать удовлетворительным, то целесообразно прибегнуть к уточнению элементов обратной матрицы по следующему правилу.

Пусть D_0 — матрица, полученная из данной матрицы A каким-либо процессом обращения. Рассмотрим матрицу R_0 , являющуюся погрешностью обращения и определяемую формулой

$$R_0 = E - AD_0, \quad (2.3.41)$$

и предположим, что $\|R_0\| \leq K < 1$. При этом условии элементы обратной матрицы A^{-1} могут быть вычислены при помощи следующего итерационного процесса со сколь угодно большой точностью. Образует две последовательности матриц $\{D_k\}$ и $\{R_k\}$, члены которых определяются по формулам

$$D_k = D_{k-1}(E + R_{k-1}), \quad R_k = E - AD_k \quad (k=1, 2, \dots). \quad (2.3.42)$$

Покажем, что последовательность матриц $\{R_k\}$ быстро убывающая и что $\lim_{m \rightarrow \infty} R_m = 0$. В силу формул (2.3.42) имеем

$$\begin{aligned} R_m &= E - AD_m = E - AD_{m-1}(E + R_{m-1}) = (E - AD_{m-1}) - AD_{m-1}R_{m-1} = \\ &= R_{m-1} - AD_{m-1}R_{m-1} = R_{m-1}^2. \end{aligned}$$

Значит,

$$\left. \begin{aligned} R_1 &= R_0^2, \\ R_2 &= R_1^2 = R_0^{2^2}, \\ &\dots \dots \dots \\ R_m &= R_{m-1}^2 = \dots = R_0^{2^m}. \end{aligned} \right\} \quad (2.3.43)$$

Поэтому $\|R_m\| \leq \|R_0\|^{2^m} \leq K^{2^m}$ и $\lim_{m \rightarrow \infty} \|R_m\| = 0$, откуда следует, что $\lim_{m \rightarrow \infty} R_m = 0$. Далее, учитывая формулы (2.3.42), (2.3.43), получим

$$D_m = A^{-1}(E - R_0^{2^m}). \quad (2.3.44)$$

Эта формула показывает, что D_m стремится при возрастании m к A^{-1} , причем сходимость процесса очень быстрая. Дадим оценку погрешности, т. е. разности $D_m - A^{-1}$. Имеем

$$D_m - A^{-1} = -A^{-1}R_0^{2^m}.$$

Так как $A^{-1} = D_0(AD_0)^{-1}$ и $AD_0 = E - R_0$, то окончательно получим

$$\begin{aligned} \|D_m - A^{-1}\| &= \|-A^{-1}R_0^{2^m}\| = \|-D_0(E - R_0)^{-1}R_0^{2^m}\| \leq \\ &\leq \|D_0\| \cdot \|(E - R_0)^{-1}\| \cdot \|R_0^{2^m}\| \leq \|D_0\| \frac{K^{2^m}}{1 - K}. \end{aligned} \quad (2.3.45)$$

Из оценки (2.3.45) видно, что если $K \ll 1$, то сходимость D_m к A^{-1} будет очень быстрой и для уточнения обратной матрицы не придется выполнять большое число итераций. Отметим еще, что члены последовательности матриц $\{D_k\}$ можно вычислять по несколько видоизмененной формуле (2.3.42), а именно:

$$D_k = D_{k-1}(E + R_{k-1}) = D_{k-1} + D_{k-1}(E - AD_{k-1}).$$

В этой формуле второе слагаемое будет играть роль поправочного члена.

Общие методы уточнения полученных решений и способы ускорения сходимости итерационных процессов при решении задач линейной алгебры будут рассмотрены нами в следующей главе.

Заканчивая параграф, обратим внимание на следующее обстоятельство. Как правило, задачи решения системы уравнений, вычисления определителей, обращения матриц будут решаться тем точнее, чем меньшей будет суммарная ошибка, вносимая при выполнении ряда операций, связанных с сильным накоплением погрешностей округлений. Это в первую очередь относится к операциям вычисления скалярного произведения, произведения матриц, произведения матрицы на вектор и т. д. Поэтому

в суммах вида $\sum_{k=1}^n \alpha_k \beta_k$ должны выполняться арифметические операции с двойной точностью и округляться должно не каждое слагаемое, а весь результат.

§ 2.4. МЕТОДЫ, ОСНОВАННЫЕ НА РАЗЛОЖЕНИЯХ МАТРИЦЫ

В предыдущем параграфе мы рассмотрели несколько методов исключения, которые в силу формулы (2.3.2) можно было трактовать также и как методы, основанные на разложениях матрицы системы (2.3.1) в произведение двух или более матриц специального вида. Однако во всех этих методах разложение матрицы в произведение матриц в явном виде не выписывается и сам вид матриц произведения по существу остается нам неизвестным и никак в схемах методов не используется. Там в полной мере присутствует только идея исключения неизвестных с помощью линейного комбинирования строк матрицы. Поэтому мы и объединили эти методы под общим названием методов исключения.

Ниже мы рассмотрим методы решения систем вида (2.3.1), в основу которых положена идея разложения искомой матрицы в произведении двух или более матриц специального вида, причем здесь в схемах методов существенную роль будут играть сами матрицы произведения, их вид, структура. Из числа таких матриц чаще всего используются матрицы особого вида, предназначенные для исключения одного или нескольких неизвестных, так называемые матрицы вращения и отражения.

2.4.1. Метод квадратного корня

Этот метод применяется при решении систем вида (2.3.1) с неособенной эрмитовой матрицей. Если матрица A не является эрмитовой, то без предварительного преобразования системы к виду $A^*A\bar{x}=A^*\bar{f}$ метод применять нельзя. Однако преобразование системы к указанному выше виду связано с выполнением большого числа дополнительных операций умножения и сложения, число которых намного превосходит число аналогичных операций, необходимых при решении системы с эрмитовой матрицей по методу квадратного корня. Поэтому выполнять указанное преобразование и затем применять к решению системы метод квадратного корня, как правило, не целесообразно.

Пусть матрица A системы $A\bar{x}=\bar{f}$ эрмитова. Схема метода квадратного корня строится на идее представления матрицы A в виде произведения треугольных и диагональных матриц, а именно: находим такую правую треугольную матрицу S и диагональную матрицу D с элементами ± 1 по главной диагонали, чтобы имело место равенство

$$A=S^*DS, \quad (2.4.1)$$

где приняты обозначения

$$S=\begin{bmatrix} s_{11} & s_{12} & \dots & s_{1n} \\ 0 & s_{22} & \dots & s_{2n} \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & s_{nn} \end{bmatrix}, \quad D=\begin{bmatrix} d_{11} & 0 & \dots & 0 \\ 0 & d_{22} & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & d_{nn} \end{bmatrix}.$$

Предположим, что мы нашли такие матрицы S и D , для которых имеет место равенство (2.4.1). Тогда решение системы $A\bar{x}=\bar{f}$ осуществляется по такому правилу. Введем следующие обозначения:

$$B=S^*D, \quad B=(\beta_{ij}), \quad S\bar{x}=\bar{y}, \quad \bar{y}=(y_1, y_2, \dots, y_n)',$$

где B — известная матрица; \bar{y} — неизвестный вектор.

Для определения \bar{y} , в силу формул

$$A\bar{x} = S^* D S \bar{x} = (S^* D) S \bar{x} = \bar{f},$$

имеем такую систему линейных алгебраических уравнений:

$$B\bar{y} = \bar{f}. \quad (2.4.2)$$

Здесь особенно важно то, что матрица этой системы является левой треугольной, т. е. имеет вид

$$B = \begin{bmatrix} \beta_{11} & 0 & & 0 \\ \beta_{21} & \beta_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{n1} & \beta_{n2} & & \beta_{nn} \end{bmatrix}.$$

Это позволяет сразу из системы (2.4.2) выписать ее решение, выполняя только обратный ход метода Гаусса сверху вниз. В результате получим

$$\left. \begin{aligned} y_1 &= \frac{f_1}{\beta_{11}}, \\ y_k &= \frac{f_k - \sum_{s=1}^{k-1} \beta_{ks} y_s}{\beta_{kk}} \quad (k=2, 3, \dots, n). \end{aligned} \right\} \quad (2.4.3)$$

Определив вектор \bar{y} , находим затем из системы $S\bar{x} = \bar{y}$ искомое решение системы $A\bar{x} = \bar{f}$. Для этого нам надо будет в системе $S\bar{x} = \bar{y}$ выполнить обратный ход метода Гаусса снизу вверх, после чего получим

$$\left. \begin{aligned} x_n &= \frac{y_n}{s_{nn}}, \\ x_k &= \frac{y_k - \sum_{p=k+1}^n s_{kp} x_p}{s_{kk}} \quad (k=n-1, n-2, \dots, 1). \end{aligned} \right\} \quad (2.4.4)$$

Как мы видим, для вычисления векторов \bar{y} и \bar{x} требуются простые, негромоздкие вычисления. Теперь, чтобы придать схеме метода окончательный вид, надо указать правило, по которому следует вычислять элементы матриц S и D . Соотношение (2.4.1) можно рассматривать как систему алгебраических уравнений для определения $\frac{n(n+1)}{2}$ элементов

матрицы S и n элементов матрицы D . Так как матрица A эрмитова, то мы будем располагать $\frac{n(n+1)}{2}$ уравнениями следующего вида:

$$\left. \begin{aligned} \bar{s}_{i1}d_{11}s_{1j} + \bar{s}_{2i}d_{22}s_{2j} + \dots + \bar{s}_{ii}d_{ii}s_{ij} &= a_{ij} \quad (i < j), \\ |s_{i1}|^2d_{11} + |s_{i2}|^2d_{22} + \dots + |s_{ii}|^2d_{ii} &= a_{ii} \quad (i = j) \end{aligned} \right\} \quad (2.4.5)$$

$$(j = 1, 2, \dots, n).$$

Здесь \bar{s}_{ij} — число, комплексно сопряженное с s_{ij} . В системе (2.4.5) число уравнений меньше числа неизвестных на n . Чтобы разложение (2.4.1) было однозначным, определим диагональные элементы s_{ii} так, чтобы они были вещественны и положительны. Тогда из второго уравнения системы (2.4.5) при $i=1$ имеем

$$|s_{11}|^2d_{11} = a_{11}.$$

Положим $d_{11} = \text{sign } a_{11}$ и из предыдущего уравнения для s_{11} получим $s_{11} = \sqrt{|a_{11}|}$. Из первого уравнения системы (2.4.5) при $i=1$ найдем $s_{1j} = \frac{a_{1j}}{d_{11}s_{11}}$ ($j=2, 3, \dots, n$). Таким образом, мы сможем определить элементы первой строки матрицы S . Далее, аналогично, из второго уравнения системы (2.4.5) и из первого уравнения при $i=2$ находим:

$$d_{22} = \text{sign}(a_{22} - |s_{12}|^2d_{11}), \quad s_{22} = \sqrt{|a_{22} - |s_{12}|^2d_{11}|},$$

$$s_{2j} = \frac{a_{2j} - \bar{s}_{12}d_{11}s_{1j}}{d_{22}s_{22}} \quad (j=3, 4, \dots, n).$$

Эти формулы позволяют вычислить элементы второй строки матрицы S . Продолжая этот процесс, мы сможем вычислить все элементы матрицы S . Укажем в общем виде формулы, по которым должны вестись вычисления элементов s_{ij} :

$$\left. \begin{aligned} d_{11} &= \text{sign } a_{11}, \quad s_{11} = \sqrt{|a_{11}|}, \quad s_{1j} = \frac{a_{1j}}{d_{11}s_{11}}, \\ d_{ii} &= \text{sign} \left(a_{ii} - \sum_{p=1}^{i-1} |s_{pi}|^2d_{pp} \right), \\ s_{ii} &= \sqrt{\left| a_{ii} - \sum_{p=1}^{i-1} |s_{pi}|^2d_{pp} \right|} \quad (i > 1), \\ s_{ij} &= \frac{a_{ij} - \sum_{p=1}^{i-1} \bar{s}_{pi}d_{pp}s_{pj}}{d_{ii}s_{ii}} \quad (i < j, j = i+1, i+2, \dots, n). \end{aligned} \right\} \quad (2.4.6)$$

Таким образом, при решении системы $A\bar{x}=\bar{f}$ по методу квадратного корня необходимо:

1) сначала убедиться в том, что A — эрмитова матрица, и затем по формулам (2.4.6) вычислить элементы матрицы S ;

2) используя формулы (2.4.3), вычислить вектор \bar{y} ;

3) наконец, по формулам (2.4.4) найти искомое решение системы $A\bar{x}=\bar{f}$ — вектор \bar{x} .

Если матрица A — вещественная и симметрическая, то ее можно разложить в произведение двух транспонированных друг другу треугольных матриц, а именно:

$$A=S'S,$$

где S — правая треугольная. В этом случае формулы (2.4.6) несколько упростятся и будут иметь вид

$$\left. \begin{aligned} s_{11} &= \sqrt{a_{11}}, \quad s_{1j} = \frac{a_{1j}}{s_{11}}, \\ s_{ii} &= \sqrt{a_{ii} - \sum_{p=1}^{i-1} s_{pi}^2} \quad (i > 1), \\ s_{ij} &= \frac{a_{ij} - \sum_{p=1}^{i-1} s_{pi}s_{pj}}{s_{ii}} \quad (i < j, j = i+1, i+2, \dots, n). \end{aligned} \right\} \quad (2.4.7)$$

Заметим, что в указанном разложении диагональные элементы матрицы S будут вещественными и положительными только в случае, когда матрица A положительно определенная. В противном случае среди элементов s_{ii} , равно как и среди других элементов s_{ij} матрицы S , могут быть и комплексные.

Для решения системы линейных алгебраических уравнений с вещественной симметрической матрицей порядка n по методу квадратного корня необходимо выполнить: умножений и делений $\frac{1}{6}n(n^2+9n+8)$, извлечений квадратных корней n .

Отметим в заключение, что метод квадратного корня очень эффективен при решении систем с положительно определенной эрмитовой матрицей. Такие системы, как правило, возникают при решении задач минимизации положительно определенных квадратичных форм. Кроме того, в методе квадратного корня имеется возможность полнее использовать другие специфические свойства матрицы A . Так, например, если матрица A имеет вид такой, как на рис. 2.4.1, a и b , то матрицы S будут иметь соответственно вид (рис. 2.4.2, a и b). Действительно, если при некотором j коэффициенты эрмитовой матрицы A удовлетворяют условию

$a_{ij}=0$ для всех $1 \leq i \leq m_j < j$, то тогда, как это следует из формул (2.4.6), и все соответствующие элементы $s_{ij}=0$. Исключение операций для этих нулевых элементов матрицы S позволяют не только решать системы быстрее, но и увеличивать порядок решаемых систем.

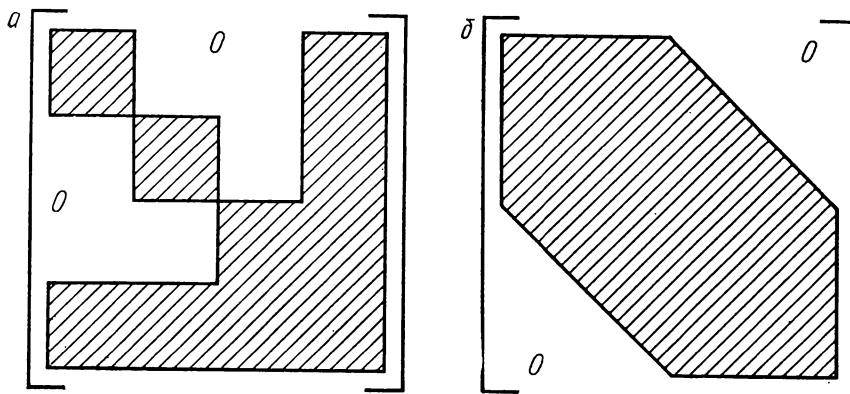


Рис. 2.4.1

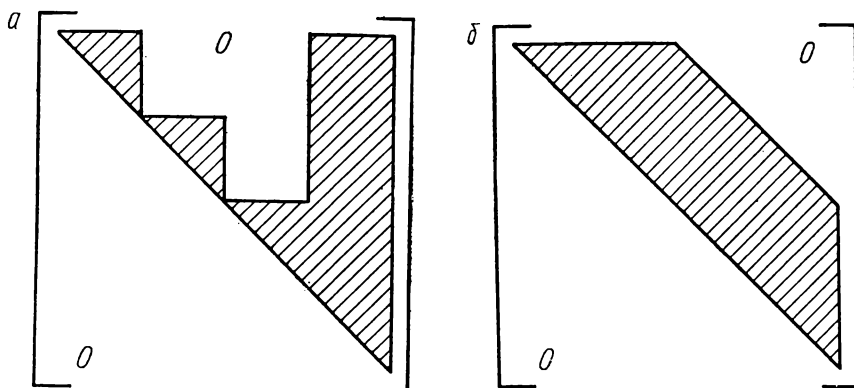


Рис. 2.4.2

В общем случае, когда матрица A системы $A\bar{x}=\bar{f}$ не является эрмитовой, к решению системы также может быть применена идея разложения матрицы A в произведение двух матриц специального вида. Основанием для этого служит следующая

Теорема 1. *Какова бы ни была матрица A с отличными от нуля главными минорами*

$$a_{11} \neq 0, \quad \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} \neq 0, \dots, \quad \begin{vmatrix} a_{11} & \dots & a_{1\ n-1} \\ \cdot & \cdot & \cdot \\ a_{n-1\ 1} & \dots & a_{n-1\ n-1} \end{vmatrix} \neq 0,$$

ее всегда можно разложить в произведение двух треугольных матриц

$$A = BC, \quad (2.4.8)$$

где B — левая треугольная матрица:

$$B = \begin{bmatrix} \beta_{11} & 0 & \dots & 0 \\ \beta_{21} & \beta_{22} & \dots & 0 \\ \cdot & \cdot & \cdot & \cdot \\ \beta_{n1} & \beta_{n2} & & \beta_{nn} \end{bmatrix},$$

C — правая треугольная матрица:

$$C = \begin{bmatrix} \gamma_{11} & \gamma_{12} & \dots & \gamma_{1n} \\ 0 & \gamma_{22} & \dots & \gamma_{2n} \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & & \gamma_{nn} \end{bmatrix}.$$

Доказательство. Отметим прежде всего, что если разложение вида (2.4.8) существует, то оно заведомо неединственно. В самом деле, если имеет место $A = BC$, то и

$$A = (BD^{-1})(DC),$$

где D — невырожденная диагональная матрица, опять будет произведением левой треугольной матрицы BD^{-1} на правую треугольную матрицу DC .

Перейдем теперь непосредственно к доказательству теоремы. Выполним в формуле (2.4.8) умножение матриц, получим

$$\sum_{k=1}^{\min(i, j)} \beta_{ik} \gamma_{kj} = a_{ij}. \quad (2.4.9)$$

Отсюда при $i=j=1$ имеем: $\beta_{11}\gamma_{11}=a_{11}$. Это уравнение позволяет определить β_{11} и γ_{11} с точностью до некоторого произвольного постоянного множителя, например, можно положить $\gamma_{11}=\delta_1 \neq 0$ и $\beta_{11}=\frac{a_{11}}{\delta_1}$. Далее при $i=j>1$ имеем

$$\beta_{ii}\gamma_{ii} + \sum_{k=1}^{i-1} \beta_{ik}\gamma_{ki} = a_{ii}. \quad (2.4.10)$$

Из (2.4.10) при $i=2$ аналогично определим β_{22} и γ_{22} . Когда $i > j$, то (2.4.9) дает следующую формулу для определения β_{ij} :

$$\beta_{ij} = \frac{a_{ij} - \sum_{k=1}^{j-1} \beta_{ik}\gamma_{kj}}{\gamma_{jj}} \quad (i > j), \quad (2.4.11)$$

а при $j > i$ (2.4.9) дает формулу для определения γ_{ij} :

$$\gamma_{ij} = \frac{a_{ij} - \sum_{k=1}^{i-1} \beta_{ik}\gamma_{kj}}{\beta_{ii}} \quad (j > i). \quad (2.4.12)$$

После определения β_{22} и γ_{22} из (2.4.11) при $i=2, j=1$ находим β_{21} , а из (2.4.12) при $i=1, j=2$ находим γ_{12} . Последовательно используя формулы (2.4.10) — (2.4.12), мы сможем определить элементы третьей строки матрицы B и третьего столбца матрицы C и т. д. Этот процесс можно осуществить только в том случае, когда $\beta_{ii}\gamma_{ii} \neq 0$ при всех $i=1, 2, \dots, n-1$, так как на эти элементы выполняется деление в формулах (2.4.11), (2.4.12). Покажем, что если матрица A имеет отличные от нуля главные миноры, то $\beta_{ii}\gamma_{ii} \neq 0$ при всех $i=1, 2, \dots, n-1$. Действительно, если бы $\beta_{ii}\gamma_{ii} = 0$ при $i=1, 2, \dots, k-1$, а $\beta_{kk}\gamma_{kk} = 0$ при некотором $k \leq n-1$, то, в силу формул (2.3.10) — (2.3.12), было бы возможно разложение

$$A_k = B_k C_k,$$

где A_k, B_k, C_k — главные миноры порядка k соответственно матриц A, B, C . Вычисляя теперь определитель матрицы A_k , мы получили бы

$$|A_k| = |B_k| \cdot |C_k| = \prod_{i=1}^k \beta_{ii}\gamma_{ii} = 0 \quad (\text{ибо } \beta_{kk}\gamma_{kk} = 0),$$

что невозможно, так как по условию теоремы $|A_k| \neq 0$. Значит, наше предположение о том, что $\beta_{kk}\gamma_{kk} = 0$ при $k \leq n-1$ не верно и разложение вида (2.3.8) существует.

Теорема доказана.

Отметим, что при фиксировании элементов по главной диагонали у матриц B или C разложение вида (2.4.8) будет единственным. Можно, например, полагать $\gamma_{11} = \gamma_{22} = \dots = \gamma_{nn} = 1$.

После того как матрица A системы $A\bar{x} = \bar{f}$ разложена в произведение вида (2.4.8), искомый вектор \bar{x} может быть вычислен по формулам типа (2.4.3), (2.4.4).

2.4.2. Метод отражений

Этот метод основан на разложении матрицы A системы $A\bar{x}=\bar{f}$ в произведение унитарной матрицы на правую треугольную матрицу. Причем здесь унитарная матрица образуется как произведение нескольких квадратных матриц, так называемых *матриц отражения*. Это название матрицы получили из-за их свойства осуществлять преобразование векторного пространства по правилу отражения векторов от заданной плоскости. Изложим правило построения таких матриц.

Пусть Q — некоторая заданная плоскость. Рассмотрим произвольный вектор

$$\bar{z}_0 = \bar{x} + \bar{y},$$

где вектор \bar{x} обладает свойством $(\bar{x}, \bar{\omega}) = 0$, $\bar{\omega}$ — вектор-столбец единичной длины, ортогональный Q , $\bar{y} = \alpha \bar{\omega}$, α — произвольное число. Вектор \bar{z}_1 , полученный как результат отражения \bar{z}_0 от Q , очевидно, имеет такой вид: $\bar{z}_1 = \bar{x} - \bar{y}$. Матрицу отражения, переводящую \bar{z}_0 в \bar{z}_1 , обозначим через V , $V\bar{z}_0 = \bar{z}_1$. Вид этой матрицы определяется формулой

$$V = E - 2\bar{\omega}\bar{\omega}^*. \quad (2.4.13)$$

Проверим, что $V\bar{z}_0 = \bar{z}_1$:

$$\begin{aligned} V\bar{z}_0 &= (E - 2\bar{\omega}\bar{\omega}^*)\bar{z}_0 = \bar{z}_0 - 2\bar{\omega}\bar{\omega}^*(\bar{x} + \alpha\bar{\omega}) = \bar{z}_0 - 2\bar{\omega}\bar{\omega}^*\bar{x} - \alpha \cdot 2\bar{\omega}\bar{\omega}^*\bar{\omega} = \\ &= \bar{z}_0 - 2\alpha\bar{\omega} = \bar{x} + \alpha\bar{\omega} - 2\alpha\bar{\omega} = \bar{x} - \alpha\bar{\omega} = \bar{z}_1, \end{aligned}$$

ибо $\bar{\omega}\bar{\omega}^*\bar{x} = \bar{\omega}(\bar{x}, \bar{\omega}) = 0$ и $\bar{\omega}\bar{\omega}^*\bar{\omega} = \bar{\omega}(\bar{\omega}, \bar{\omega}) = \bar{\omega}$. Легко проверяется и тот факт, что V — унитарная матрица:

$$\begin{aligned} VV^* &= (E - 2\bar{\omega}\bar{\omega}^*)(E - 2\bar{\omega}\bar{\omega}^*) = E - 2\bar{\omega}\bar{\omega}^* - 2\bar{\omega}\bar{\omega}^* + 4\bar{\omega}\bar{\omega}^*\bar{\omega}\bar{\omega}^* = \\ &= E - 4\bar{\omega}\bar{\omega}^* + 4\bar{\omega}(\bar{\omega}^*\bar{\omega})\bar{\omega}^* = E - 4\bar{\omega}\bar{\omega}^* + 4\bar{\omega}(\bar{\omega}^*) = E. \end{aligned}$$

Матрицы отражения V могут быть эффективно использованы при решении задачи о приведении заданной матрицы к виду правой треугольной. Чтобы показать это, рассмотрим сначала, как с помощью матрицы V перевести произвольный вектор \bar{s} в заданный вектор \bar{l} единичной длины, т. е. как определить матрицу V и число α , чтобы имело место равенство,

$$V\bar{s} = \bar{l}. \quad (2.4.14)$$

Его можно также записать в таком виде:

$$2(\bar{s}, \bar{\omega})\bar{\omega} = \bar{s} - \alpha\bar{l} \quad (2.4.15)$$

или

$$\bar{\omega} = \kappa(\bar{s} - \alpha\bar{l}),$$

где $\kappa = \frac{1}{2(\bar{s}, \bar{\omega})}$. Подставив это выражение для $\bar{\omega}$ в формулу (2.4.15), получим:

$$2(\bar{s}, \kappa(\bar{s} - \alpha\bar{l})) \cdot \kappa(\bar{s} - \alpha\bar{l}) = \bar{s} - \alpha\bar{l}$$

или

$$[2|\kappa|^2(\bar{s}, \bar{s} - \alpha\bar{l}) - 1](\bar{s} - \alpha\bar{l}) = 0.$$

Выберем κ таким образом, чтобы выражение в квадратной скобке обратилось в нуль. Это дает

$$|\kappa|^2 = \frac{1}{2(\bar{s}, \bar{s} - \alpha \bar{l})}.$$

Здесь число α подлежит выбору. Определим его таким образом, чтобы $(\bar{s}, \bar{s} - \alpha \bar{l}) > 0$. Имеем $(\bar{s}, \bar{s} - \alpha \bar{l}) = (\bar{s}, \bar{s}) - \alpha(\bar{s}, \bar{l})$. Положим $|\alpha| = \sqrt{(\bar{s}, \bar{s})}$. Тогда

$$\begin{aligned} (\bar{s}, \bar{s} - \alpha \bar{l}) &= |\alpha|^2 - \alpha(\bar{s}, \bar{l}) = |\alpha|^2 - |\alpha| e^{-i \arg \alpha} |(\bar{s}, \bar{l})| e^{i \arg(\bar{s}, \bar{l})} = \\ &= |\alpha|^2 - |\alpha| |(\bar{s}, \bar{l})| e^{i(-\arg \alpha + \arg(\bar{s}, \bar{l}))}. \end{aligned}$$

Отсюда следует, что $(\bar{s}, \bar{s} - \alpha \bar{l})$ заведомо будет больше нуля, если

$$-e^{i(-\arg \alpha + \arg(\bar{s}, \bar{l}))} = 1.$$

Для этого достаточно взять $-\arg \alpha + \arg(\bar{s}, \bar{l}) = \pi$, т. е.

$$\arg \alpha = \pi - \arg(\bar{s}, \bar{l}).$$

Тогда окончательно получим:

$$(\bar{s}, \bar{s} - \alpha \bar{l}) = |\alpha|^2 + |\alpha| |(\bar{s}, \bar{l})|$$

и

$$|\kappa|^2 = \frac{1}{2[|\alpha|^2 + |\alpha| |(\bar{s}, \bar{l})|]}.$$

Таким образом, для того чтобы матрица $V = E - 2\bar{\omega} \bar{\omega}^*$ удовлетворяла условию $\bar{V}\bar{s} = \alpha \bar{l}$, где \bar{s} и \bar{l} — заданные векторы, надо положить

$$\bar{\omega} = \kappa(\bar{s} - \alpha \bar{l}), \quad \alpha = \sqrt{(\bar{s}, \bar{s})}, \quad \kappa = \frac{1}{\sqrt{2(\bar{s} - \alpha \bar{l}, \bar{s} - \alpha \bar{l})}} = \frac{1}{\sqrt{2(|\alpha|^2 + |\alpha| |(\bar{s}, \bar{l})|)}}.$$

Теперь задача разложения произвольной неособенной комплексной матрицы A в произведение унитарной и правой треугольной матриц решается так. На первом шаге образуем матрицу V_1 , взяв в качестве \bar{s} и \bar{l} следующие векторы: $\bar{s} = (a_{11}, a_{21}, \dots, a_{n1})'$, $\bar{l} = (1, 0, \dots, 0)'$ и вычислив $\bar{\omega}$, α и κ по указанным выше формулам. Умножив A слева на V_1 , мы придем к матрице $A^{(1)}$ вида

$$A^{(1)} = V_1 A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & & a_{1n}^{(1)} \\ 0 & a_{22}^{(1)} & \dots & a_{2n}^{(1)} \\ \dots & \dots & \dots & \dots \\ 0 & a_{n2}^{(1)} & \dots & a_{nn}^{(1)} \end{bmatrix}.$$

Очевидно, что $a_{11}^{(1)} = \alpha$. На втором шаге аналогичным путем образуем матрицу V_2 по векторам $\bar{s} = (0, a_{22}^{(1)}, \dots, a_{n2}^{(1)})'$, $\bar{l} = (0, 1, 0, \dots, 0)'$ и выполним умножение слева $A^{(1)}$ на V_2 , в результате чего получим матрицу $A^{(2)}$ вида

$$A^{(2)} = V_2 A^{(1)} = V_2 V_1 A = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} & a_{13}^{(2)} & \dots & a_{1n}^{(2)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \dots & a_{2n}^{(2)} \\ 0 & 0 & a_{33}^{(2)} & \dots & a_{3n}^{(2)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & a_{n3}^{(2)} & \dots & a_{nn}^{(2)} \end{bmatrix}.$$

В этой матрице первая строка имеет тот же вид, что и аналогичная строка в матрице $A^{(1)}$, в силу того, что матрица V_2 имеет такой вид:

$$V_2 = E - 2\bar{\omega}\bar{\omega}^* = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & u_{22}^{(2)} & u_{23}^{(2)} & \dots & u_{2n}^{(2)} \\ 0 & u_{32}^{(2)} & u_{33}^{(2)} & \dots & u_{3n}^{(2)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & u_{n2}^{(2)} & u_{n3}^{(2)} & \dots & u_{nn}^{(2)} \end{bmatrix}.$$

Продолжая этот процесс дальше, мы на n -м шаге придем к матрице $A^{(n)}$ вида

$$A^{(n)} = V_n V_{n-1} \dots V_2 V_1 A = \begin{bmatrix} a_{11}^{(n)} & a_{12}^{(n)} & \dots & a_{1n-1}^{(n)} & a_{1n}^{(n)} \\ 0 & a_{22}^{(n)} & \dots & a_{2n-1}^{(n)} & a_{2n}^{(n)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0 & a_{nn}^{(n)} \end{bmatrix}.$$

Обозначив здесь $V_n V_{n-1} \dots V_1$ через V , мы получим

$$A^{(n)} = VA,$$

откуда уже следует искомое представление матрицы A в виде произведения унитарной матрицы на правую треугольную матрицу, т. е.

$$A = V^* A^{(n)}.$$

Основываясь на изложенной выше теории, построим вычислительную схему метода отражений. Пусть требуется решить систему

$$A\bar{x} = \bar{f}$$

с неособенной комплексной матрицей A . Рассмотрим расширенную матрицу этой системы со столбцами $\bar{a}_1^{(0)}, \bar{a}_2^{(0)}, \dots, \bar{a}_n^{(0)}, \bar{a}_{n+1}^{(0)}$ и обозначим ее через A_0 , таким образом,

$$A_0 = [\bar{a}_1^{(0)}, \bar{a}_2^{(0)}; \dots, \bar{a}_{n+1}^{(0)}],$$

где

$$\bar{a}_k^{(0)} = (a_{1k}, a_{2k}, \dots, a_{nk})' \quad (k=1, 2, \dots, n), \quad \bar{a}_{n+1}^{(0)} = \bar{f}.$$

Отметим, что невозможность выполнения очередного шага связана лишь с равенством нулю очередного вектора \bar{s} , ибо тогда $\alpha=0$ и нельзя вычислить число ω . Однако вектор \bar{s} не может быть нулевым, так как матрица A невырождена и преобразуется она унитарными матрицами. Рекомендуется с целью уменьшения общего объема вычислений формулу (2.4.17) использовать в такой форме:

$$\bar{a}_i^{(k+1)} = \bar{a}_i^{(k)} - 2(\bar{a}_i^{(k)}, \bar{\omega}) \bar{\omega}.$$

Этот вид формулы следует из (2.4.17), если учесть, что $(\bar{\omega} \bar{\omega}^*) \bar{a}_i^{(k)} = (\bar{a}_i^{(k)}, \bar{\omega}) \bar{\omega}$.

Для решения системы линейных алгебраических уравнений методом отражений необходимо выполнить $\frac{1}{6}(4n^3 + 15n^2 + 23n - 6)$ умножений и делений и $n-1$ извлечений квадратных корней.

В изложенном выше методе отражений исключение неизвестных на каждом шаге производилось с помощью матриц отражения. Обратим внимание на то, что эта же задача может быть решена и с помощью элементарных унитарных матриц $R_{ij}(\varphi, \psi)$ вида

$$R_{ij}(\varphi, \psi) = \begin{bmatrix} \overset{(i)}{1} & & & & & & & \overset{(j)}{0} \\ & \ddots & & & & & & \\ & & \cos \varphi & & -e^{i\psi} \sin \varphi & & & \\ & & & \ddots & & & & \\ & & & & 1 & & & \\ & & & & & \ddots & & \\ & & & & & & 1 & \\ & & e^{-i\psi} \sin \varphi & & & \cos \varphi & & \\ & & & & & & \ddots & \\ 0 & & & & & & & 1 & \ddots & 1 \end{bmatrix} \begin{matrix} \\ \\ (i) \\ \\ \\ \\ (j) \\ \\ \end{matrix}.$$

Действительно, если A — произвольная неособенная комплексная матрица, то, умножив ее слева на матрицу $R_{ij}(\varphi, \psi)$, мы получим новую матрицу B , у которой элементы i -й и j -й строк определяются по формулам

$$\left. \begin{aligned} b_{ip} &= \cos \varphi \cdot a_{ip} - \sin \varphi e^{i\psi} \cdot a_{jp}, \\ b_{jp} &= \sin \varphi \cdot e^{-i\psi} a_{ip} + \cos \varphi \cdot a_{jp} \quad (p=1, 2, \dots, n), \end{aligned} \right\} \quad (2.4.19)$$

а остальные элементы матрицы B такие же, как и у матрицы A . Если теперь мы хотим элемент b_{js} матрицы B обратить в нуль (это эквивалентно исключению из j -го уравнения неизвестного x_s с помощью операции $R_{ij}A\bar{x} = R_{ij}\bar{f}$), то необходимо в формуле (2.4.19) при $p=s$ взять

$$\left. \begin{aligned} \psi &= \arg a_{is} - \arg a_{js}, \\ \cos \varphi &= \frac{|a_{is}|}{\sqrt{|a_{is}|^2 + |a_{js}|^2}}, \\ \sin \varphi &= \frac{-|a_{js}|}{\sqrt{|a_{is}|^2 + |a_{js}|^2}}, \end{aligned} \right\} \quad (2.4.20)$$

если $\sqrt{|a_{is}|^2 + |a_{js}|^2} \neq 0$. В противном случае следует положить

$$\cos \varphi = 1, \quad \sin \varphi = 0.$$

Такое свойство матриц $R_{ij}(\varphi, \psi)$ позволяет утверждать, что справедлива

Теорема 2. Любая комплексная матрица A преобразуется в правую треугольную матрицу посредством умножения слева на конечную цепочку матриц $R_{ij}(\varphi, \psi)$.

Доказательство. Умножим матрицу A слева на матрицы $R_{12}, R_{13}, \dots, R_{1n}$, выбирая их так, чтобы последовательно аннулировать все поддиагональные элементы первого столбца. В результате мы получим:

$$A^{(1)} = R_{1n} R_{1n-1} \dots R_{12} A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(1)} & \dots & a_{2n}^{(1)} \\ \cdot & \cdot & \cdot & \cdot \\ 0 & a_{n2}^{(1)} & \dots & a_{nn}^{(1)} \end{bmatrix}.$$

На втором шаге матрицу $A^{(1)}$ умножаем на соответствующим образом подобранные матрицы $R_{23}, R_{24}, \dots, R_{2n}$, на третьем шаге — на $R_{33}, R_{34}, \dots, R_{3n}$ и т. д. В конце процесса мы получим правую треугольную матрицу $A^{(n-1)}$ по формуле

$$A^{(n-1)} = R_{n-1n} R_{n-2n} R_{n-2n-1} \dots R_{12} A. \quad (2.4.21)$$

Теорема доказана.

Из формулы (2.4.21) следует уже установленный нами ранее факт, что любая комплексная матрица есть произведение унитарной на правую треугольную. Действительно, матрица $V = R_{n-1n} R_{n-2n} R_{n-2n-1} \dots R_{12}$ — унитарна, поэтому из (2.4.21) получим:

$$A = V^* A^{(n-1)},$$

но V^* также унитарна.

Сравнивая описанный процесс приведения матрицы к виду правой треугольной с аналогичным процессом, осуществляемым с помощью матриц отражения, мы видим, что применение матриц отражения в этой задаче оказывается более эффективным, ибо с их помощью на каждом шаге можно обращать в нуль все поддиагональные элементы некоторого столбца.

2.4.3. Вычисление определителей

Метод квадратного корня и метод отражений при решении системы уравнений могут быть попутно использованы также и для вычисления определителя матрицы. Действительно, из формулы (2.4.1) имеем

$$|A| = |S^*| \cdot |D| \cdot |S| = \prod_{i=1}^n \bar{s}_{ii} d_{ii} s_{ii}.$$

Значит,

$$|A| = \prod_{i=1}^n |s_{ii}|^2 d_{ii}. \quad (2.4.22)$$

Формула (2.4.22) является искомой для вычисления определителя матрицы A в методе квадратного корня.

Применяя формулу (2.4.16) $n-1$ раз, получим:

$$A_{n-1} = V_{n-1} V_{n-2} \dots V_1 A.$$

Отсюда

$$|A_{n-1}| = \prod_{i=1}^{n-1} |V_i| \cdot |A|.$$

Значит,

$$|A| = \frac{|A_{n-1}|}{\prod_{i=1}^{n-1} |V_i|} = \frac{\prod_{i=1}^n a_{ii}^{(n-1)}}{\prod_{i=1}^{n-1} |V_i|} \quad (2.4.23)$$

По формуле (2.4.23) можно вычислить определитель матрицы A , если мы сможем определить значение определителей матриц V_i .

Докажем, что, независимо от выбора вектора $\bar{\omega}$, определитель матрицы отражения равен минус единице, т. е. $|V_i| = -1$ при любом i . Рассмотрим матрицу $V = E - 2\bar{\omega}\bar{\omega}^*$ и покажем, что $|V| = -1$, если $\bar{\omega}$ — вектор единичной длины. Пусть $\lambda_i(V)$ — собственные числа матрицы V . Тогда

$$|V| = \prod_{i=1}^n \lambda_i(V). \quad (2.4.24)$$

Изучим свойства чисел $\lambda_i(V)$. С этой целью введем в рассмотрение матрицу $W = \bar{\omega}\bar{\omega}^*$ и собственные числа $\lambda_i(W)$ этой матрицы. Очевидно, что

$$\lambda_i(V) = 1 - 2\lambda_i(W) \quad (i = 1, 2, \dots, n).$$

Если теперь мы определим $\lambda_i(W)$, то тем самым будут определены и $\lambda_i(V)$. Матрица W — нормальная, т. е. она перестановочна со своей сопряженной. Действительно,

$$W^*W - WW^* = (\bar{\omega}\bar{\omega}^*)^*(\bar{\omega}\bar{\omega}^*) - (\bar{\omega}\bar{\omega}^*)(\bar{\omega}\bar{\omega}^*)^* = 0.$$

Для такой матрицы известна связь между квадратом ее сферической нормы и собственными значениями [2, стр. 26], а именно:

$$N^2(W) = \sum_{i=1}^n |\lambda_i(W)|^2,$$

где

$$N^2(W) = \sum_{i,j=1}^n |\omega_i \omega_j|^2 = \text{Sp } W^*W, \quad W = (\omega_i \omega_j).$$

Но

$$\begin{aligned} \text{Sp } W^*W &= \text{Sp } (\bar{\omega}\bar{\omega}^*)^*(\bar{\omega}\bar{\omega}^*) = \text{Sp } (\bar{\omega}\bar{\omega}^* \omega \omega^*) = \text{Sp } (\bar{\omega}\bar{\omega}^* \omega) \bar{\omega}^* = \\ &= \text{Sp } \bar{\omega} \bar{\omega}^* = \sum_{i=1}^n \omega_i \bar{\omega}_i = (\bar{\omega}, \bar{\omega}) = 1, \end{aligned}$$

значит,

$$N^2(W) = \text{Sp } W^*W = \sum_{i=1}^n |\lambda_i(W)|^2 = 1. \quad (2.4.25)$$

Так как матрица W имеет одно собственное значение, равное единице, ибо

$$W\bar{\omega} = \bar{\omega},$$

то в силу (2.4.25) собственные значения W распределены с точностью до нумерации следующим образом:

$$\left. \begin{aligned} \lambda_1(W) &= 1, \\ \lambda_k(W) &= 0 \quad (k=2, 3, \dots, n). \end{aligned} \right\} \quad (2.4.26)$$

Следовательно, можно считать, что

$$\left. \begin{aligned} \lambda_1(V) &= -1, \\ \lambda_k(V) &= 1 \quad (k=2, 3, \dots, n). \end{aligned} \right\}$$

Из (2.4.24) теперь окончательно получим:

$$|V| = -1.$$

Это позволяет переписать формулу (2.4.23) так:

$$|A| = (-1)^{n-1} \prod_{i=1}^n a_{ii}^{(n-1)} = - \prod_{i=1}^n (-a_{ii}^{(n-1)}). \quad (2.4.27)$$

Формула (2.4.27) является искомой для вычисления определителя в методе отражений.

Отметим, что при вычислении определителей по формулам (2.4.22), (2.4.27) следует пользоваться приемом, изложенным в п. 2.3.4, с тем чтобы избежать получения машинного нуля или переполнения.

2.4.4. Обращение матриц

Вычисления, проводимые при решении системы уравнений по методу квадратного корня, не представляется возможным непосредственно использовать в задаче обращения матрицы. Однако этот метод можно применить для нахождения обратной матрицы в случае неособенных эрмитовых матриц аналогично тому, как это мы делали в методе Гаусса (п. 2.3.5).

Более эффективно могут быть использованы для обращения матрицы вычисления, проводимые в методе отражений. Действительно, применяя формулу (2.4.16) $n-1$ раз, получим

$$A_{n-1} = V_{n-1} V_{n-2} \dots V_1 A, \quad (2.4.28)$$

где A_{n-1} — правая треугольная матрица. Ее элементы известны. Из (2.4.28) находим:

$$A^{-1} = A_{n-1}^{-1} V_{n-1} V_{n-2} \dots V_1. \quad (2.4.29)$$

Формула (2.4.29) является искомой для обращения матрицы по методу отражений. Заметим, что поскольку матрица A_{n-1} — правая треугольная, то ее следует обращать по упрощенной схеме метода окаймления, указанной в п. 2.3.5. При реализации формулы (2.4.29) нет необходимости матрицы V_k запоминать отдельно и хранить в памяти машины. Целесообразнее хранить в памяти машины только соответствующие векторы $\bar{\omega}$, а матрицы V_k вычислять по формуле (2.4.13), используя нужный вектор $\bar{\omega}$. Это позволит существенно экономить память машины, так как для хранения векторов $\bar{\omega}$ нужно всего

иметь лишь $\frac{n(n+1)}{2} - 1$ ячеек.

§ 2.5. МЕТОДЫ, ОСНОВАННЫЕ НА ПОСТРОЕНИИ ВСПОМОГАТЕЛЬНОЙ СИСТЕМЫ ВЕКТОРОВ, ОРТОГОНАЛЬНЫХ В НЕКОТОРОЙ МЕТРИКЕ

Основной особенностью рассматриваемых ниже методов является то, что в них искомое решение определяется как последний вектор в специальном образом построенной вспомогательной системе векторов. В методе ортогонализации, например, таким будет вектор, ортогональный к подпространству, натянутому на векторы-строки матрицы системы уравнений, и имеющий последнюю координату, равную единице; в методе сопряженных градиентов таким будет вектор последовательных приближений к решению системы, обращающий в нуль один из ортогональных векторов-невязок системы. Оба упомянутых метода позволяют получить точное решение системы n линейных алгебраических уравнений не позже n -го шага преобразований.

По своей идее эти методы сильно отличаются друг от друга, однако схемы их реализации имеют многие общие черты, обусловленные в основном процессом ортогонализации, проводимым в обоих методах.

В этом же параграфе мы изложим метод скорейшего спуска, который является итерационным и, следовательно, не позволяет получить точное решение за конечное число шагов, как это имело место в предыдущих методах. Целесообразность изложения этого метода здесь объясняется тем, что по своей структуре он очень тесно связан с методом сопряженных градиентов и может рассматриваться как упрощенный вариант этого метода.

2.5.1. Метод ортогонализации

Рассмотрим систему линейных алгебраических уравнений с неособенной матрицей

$$\begin{array}{rcl} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n & = & f_1, \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n & = & f_2, \\ \cdot & & \cdot \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n & = & f_n \end{array}$$

и запишем ее в таком виде:

$$\left. \begin{aligned} (\bar{a}_1, \bar{y}) &= 0, \\ (\bar{a}_2, \bar{y}) &= 0, \\ &\vdots \\ (\bar{a}_n, \bar{y}) &= 0, \end{aligned} \right\} \quad (2.5.1)$$

где приняты следующие обозначения:

$$\begin{aligned}\bar{a}_k &= (a_{k1}, a_{k2}, \dots, a_{kn}, a_{kn+1})', \quad k=1, 2, \dots, n, \quad a_{kn+1} = -f_k, \\ \bar{y} &= (x_1, x_2, \dots, x_n, 1)'. \end{aligned}$$

Из (2.5.1) следует, что решение системы линейных алгебраических уравнений с неособенной матрицей A сводится к вычислению такого вектора \bar{y} , который был бы ортогонален к линейно независимым векторам $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ и имел последнюю координату, равную единице. Ортогональность вектора \bar{y} к векторам $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ влечет за собой также ортогональность \bar{y} ко всему подпространству P_n , натянутому на векторы $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$, и, следовательно, к любому базису этого подпространства. И наоборот, ортогональность \bar{y} к любому базису P_n влечет за собой ортогональность вектора \bar{y} ко всему подпространству P_n и, следовательно, ортогональность к векторам $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$. Это обстоятельство позволяет указать следующий путь для вычисления вектора \bar{y} . Строим какой-либо ортогональный базис подпространства P_n и находим вектор \bar{z} , ортогональный к этому базису. Тогда, если $z^{(n+1)}$ — последняя координата вектора \bar{z} , то для \bar{y} получим искомую формулу

$$\bar{y} = \frac{\bar{z}}{z^{(n+1)}},$$

из которой находится решение системы уравнений в таком виде:

$$x_i = \frac{z^{(i)}}{z^{(n+1)}} \quad (i=1, 2, \dots, n).$$

Все это определяет такой способ решения рассматриваемой задачи. Добавим к системе линейно независимых векторов $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n$ еще один линейно независимый вектор. Таким вектором, как в этом легко убедиться, будет вектор \bar{a}_{n+1} вида

$$\bar{a}_{n+1} = (\underbrace{0, 0, \dots, 0}_n, 1)'.$$

Будем строить систему ортонормированных векторов $\bar{b}_1, \bar{b}_2, \dots, \bar{b}_{n+1}$ таких, что для любого k ($1 \leq k \leq n+1$) последовательность векторов $\bar{b}_1, \bar{b}_2, \dots, \bar{b}_k$ будет являться ортонормированным базисом подпространства P_k , натянутого на векторы $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_k$. В этом случае вектор \bar{b}_{n+1} будет ортогонален к пространству P_n , натянутому на векторы $\bar{a}_1, \bar{a}_2, \dots,$

\bar{a}_n , и поэтому искомое решение системы уравнений можно будет вычислять по формуле

$$x_i = \frac{b_i^{(n+1)}}{b_{n+1}^{(n+1)}} \quad (i=1, 2, \dots, n), \quad (2.5.2)$$

где $b_k^{(n+1)}$ ($k=1, 2, \dots, n+1$) — компоненты вектора \bar{b}_{n+1} .

Укажем теперь правило Шмидта для построения ортонормированного базиса пространства, натянутого на заданные линейно независимые векторы. Обозначим через $\bar{u}_1, \bar{u}_2, \dots, \bar{u}_k$ ортогональный базис P_k , а через $\bar{b}_1, \bar{b}_2, \dots, \bar{b}_k$ — ортонормированный в евклидовой метрике базис того же пространства.

При $k=1$ имеем

$$\bar{u}_1 = \bar{a}_1 \quad \text{и} \quad \bar{b}_1 = \frac{\bar{u}_1}{\|\bar{u}_1\|_{\text{III}}},$$

где $\|\bar{u}_1\|_{\text{III}} = \sqrt{(\bar{u}_1, \bar{u}_1)}$. Предположим, что для $k > 1$ мы построили ортогональный базис (векторы $\bar{u}_1, \bar{u}_2, \dots, \bar{u}_k$) и ортонормированный базис (векторы $\bar{b}_1, \bar{b}_2, \dots, \bar{b}_k$) подпространства P_k . Как теперь вычислить следующие векторы \bar{u}_{k+1} и \bar{b}_{k+1} для P_{k+1} ? Будем разыскивать вектор \bar{u}_{k+1} в виде

$$\bar{u}_{k+1} = \bar{a}_{k+1} + \sum_{i=1}^k c_i^{(k)} \bar{b}_i, \quad (2.5.3)$$

где $c_i^{(k)}$ — неизвестные величины, которые по смыслу задачи следует определить таким образом, чтобы выполнялось условие

$$(\bar{u}_{k+1}, \bar{u}_i) = 0 \quad (i=1, 2, \dots, k)$$

или, что то же самое, условие

$$(\bar{u}_{k+1}, \bar{b}_i) = 0 \quad (i=1, 2, \dots, k). \quad (2.5.4)$$

Используя (2.5.3), (2.5.4), получим

$$(\bar{a}_{k+1}, \bar{b}_s) + \sum_{i=1}^k c_i^{(k)} (\bar{b}_i, \bar{b}_s) = 0.$$

Но

$$(\bar{b}_i, \bar{b}_s) = \begin{cases} 0, & i \neq s, \\ 1, & i = s. \end{cases}$$

Значит,

$$c_s^{(k)} = -(\bar{a}_{k+1}, \bar{b}_s) \quad (s=1, 2, \dots, k). \quad (2.5.5)$$

Таким образом, для \bar{u}_{k+1} имеем

$$\bar{u}_{k+1} = \bar{a}_{k+1} - \sum_{i=1}^k (\bar{a}_{k+1}, \bar{b}_i) \bar{b}_i. \quad (2.5.6)$$

Вектор \bar{b}_{k+1} получим, нормируя \bar{u}_{k+1} :

$$\bar{b}_{k+1} = \frac{\bar{u}_{k+1}}{\|\bar{u}_{k+1}\|_{\text{III}}}. \quad (2.5.7)$$

Изложим теперь весь алгоритм метода ортогонализации:

- 1) вычисляем векторы \bar{u}_1 и \bar{b}_1 ;
- 2) по формуле (2.5.5) при $k=1$ вычисляем $c_1^{(1)}$ и по формулам (2.5.6), (2.5.7) вычисляем \bar{u}_2, \bar{b}_2 ;
- ...
- $n+1$) по формуле (2.5.5) при $k=n$ вычисляем $c_1^{(n)}, c_2^{(n)}, \dots, c_n^{(n)}$ и по формулам (2.5.6), (2.5.7) вычисляем $\bar{u}_{n+1}, \bar{b}_{n+1}$.

Зная вектор \bar{b}_{n+1} , искомое решение системы уравнений вычисляем по формуле (2.5.2).

В методе ортогонализации для нахождения решения системы n уравнений необходимо выполнить $n^3 + n^2 + n$ операций умножения и деления и n извлечений квадратного корня.

Метод ортогонализации легко реализуется на ЭВМ и это является одним из его основных достоинств. Однако удовлетворительные по точности результаты этот метод позволяет получать не для всех систем уравнений с неособенной матрицей. Основная причина этого явления лежит в неустойчивости процесса вычисления векторов \bar{u}_{k+1} по формуле (2.5.6), из-за которой нарушается основное свойство этих векторов — ортогональность. Остановимся на этом вопросе подробнее.

С целью упрощения вычислений предположим матрицу A системы вещественной и будем считать, что для некоторого $k \geq 1$ нами вычислены векторы $\bar{g}_1, \bar{g}_2, \dots, \bar{g}_k$, являющиеся приближенными значениями соответственно векторов $\bar{b}_1, \bar{b}_2, \dots, \bar{b}_k$. Предположим при этом, что

$$\bar{b}_i - \bar{g}_i = \bar{\epsilon}_i$$

и

$$\max_{1 \leq i \leq k} \|\bar{\epsilon}_i\|_{\text{III}} \leq \epsilon, \quad (2.5.8)$$

где ε — некоторое положительное число. Для векторов \bar{g}_i свойство ортонормированности нарушается из-за погрешностей $\bar{\varepsilon}_i$, поэтому для скалярного произведения (\bar{g}_i, \bar{g}_j) мы получим:

$$(\bar{g}_i, \bar{g}_j) = (\bar{b}_i - \bar{\varepsilon}_i, \bar{b}_j - \bar{\varepsilon}_j) = \delta_{ij} + \varepsilon_{ij}, \quad (2.5.9)$$

где δ_{ij} — символ Кронекера, $\varepsilon_{ij} = -(\bar{\varepsilon}_i, \bar{b}_j) - (\bar{b}_i, \bar{\varepsilon}_j) + (\bar{\varepsilon}_i, \bar{\varepsilon}_j)$ — числа, которые в силу (2.5.8) для всех $i, j \leq k$ будут величинами порядка ε . Условимся это записывать так: $|\varepsilon_{ij}| = O(\varepsilon)$.

Ясно, что система векторов $\bar{g}_1, \bar{g}_2, \dots, \bar{g}_k$ будет ортонормированной в том и только в том случае, когда все числа ε_{ij} будут равны нулю. Поскольку ошибки округлений носят случайный характер, то нельзя заранее предсказать поведение погрешностей $\bar{\varepsilon}_i$, а следовательно, и чисел ε_{ij} , нельзя также надеяться на то, что все числа ε_{ij} будут равняться нулю. В связи с этим представляют интерес поведение погрешности при вычислении следующего вектора \bar{g}_{k+1} и оценка чисел $\varepsilon_{k+1 j}$.

По определению

$$\varepsilon_{k+1 j} = (\bar{g}_{k+1}, \bar{g}_j) - \delta_{k+1 j}$$

или

$$\varepsilon_{k+1 j} = (\bar{g}_{k+1}, \bar{g}_j), \quad (2.5.10)$$

ибо $\delta_{k+1 j} = 0$ при $1 \leq j \leq k$. Чтобы дать оценку $\varepsilon_{k+1 j}$, надо указать формулу, по которой реально вычисляется вектор \bar{g}_{k+1} . В силу формул (2.5.6), (2.5.7)

$$\bar{u}_{k+1} = \bar{a}_{k+1} - \sum_{i=1}^k (\bar{a}_{k+1}, \bar{g}_i) \bar{g}_i$$

и

$$\bar{g}_{k+1} = \frac{\bar{u}_{k+1}}{\|\bar{u}_{k+1}\|_{III}}. \quad (2.5.11)$$

Из (2.5.10), (2.5.11) следует, что

$$\varepsilon_{k+1 j} = \frac{(\bar{u}_{k+1}, \bar{g}_j)}{\|\bar{u}_{k+1}\|_{III}} \quad (j = 1, 2, \dots, k). \quad (2.5.12)$$

Вычислим евклидову норму вектора \bar{u}_{k+1} :

$$\begin{aligned} \|\bar{u}_{k+1}\|_{\text{III}} &= \sqrt{(\bar{u}_{k+1}, \bar{u}_{k+1})} = \\ &= \left[(\bar{a}_{k+1} - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{g}_i) \bar{g}_i, \bar{a}_{k+1} - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{g}_i) \bar{g}_i) \right]^{\frac{1}{2}} \end{aligned}$$

Учитывая, что $\bar{g}_i = \bar{b}_i - \bar{\varepsilon}_i$, и оценку (2.5.8), получим:

$$\|\bar{u}_{k+1}\|_{\text{III}} = \left[(\bar{a}_{k+1}, \bar{a}_{k+1}) - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i)^2 \right]^{\frac{1}{2}} + O(\varepsilon). \quad (2.5.13)$$

Приступим теперь к вычислению ε_{k+1j} . Из (2.5.12), (2.5.13) имеем

$$\varepsilon_{k+1j} = \frac{(\bar{u}_{k+1}, \bar{g}_j)}{\|\bar{u}_{k+1}\|_{\text{III}}} = \frac{(\bar{a}_{k+1} - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{g}_i) \bar{g}_i, \bar{g}_j)}{[(\bar{a}_{k+1}, \bar{a}_{k+1}) - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i)^2]^{\frac{1}{2}} + O(\varepsilon)}.$$

Но

$$\begin{aligned} (\bar{a}_{k+1} - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{g}_i) \bar{g}_i, \bar{g}_j) &= (\bar{a}_{k+1}, \bar{b}_j - \bar{\varepsilon}_j) - \\ &- (\sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i - \bar{\varepsilon}_i) (\bar{b}_i - \bar{\varepsilon}_i), \bar{b}_j - \bar{\varepsilon}_j) = (\bar{a}_{k+1}, \bar{b}_j - \bar{\varepsilon}_j) - \\ &- \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i - \bar{\varepsilon}_i) (\delta_{ij} + \varepsilon_{ij}) = - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i - \bar{\varepsilon}_i) \varepsilon_{ij} = \\ &= - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i) \varepsilon_{ij} + \sum_{i=1}^h (\bar{a}_{k+1}, \bar{\varepsilon}_i) \varepsilon_{ij}. \end{aligned}$$

Окончательно, для ε_{k+1j} получим

$$\varepsilon_{k+1j} = \frac{- \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i) \varepsilon_{ij} + \sum_{i=1}^h (\bar{a}_{k+1}, \bar{\varepsilon}_i) \varepsilon_{ij}}{[(\bar{a}_{k+1}, \bar{a}_{k+1}) - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i)^2]^{\frac{1}{2}} + O(\varepsilon)}$$

или

$$\varepsilon_{k+1j} = \frac{- \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i) \varepsilon_{ij}}{[(\bar{a}_{k+1}, \bar{a}_{k+1}) - \sum_{i=1}^h (\bar{a}_{k+1}, \bar{b}_i)^2]^{\frac{1}{2}}} + O(\varepsilon^2). \quad (2.5.14)$$

Чтобы иметь суждение о поведении погрешности ε_{k+1j} , нам необходимо выяснить смысл коэффициентов, стоящих в формуле (2.5.14) при ε_{ij} , т. е. смысл чисел

$$\gamma_i = \frac{-(\bar{a}_{k+1}, \bar{b}_i)}{[(\bar{a}_{k+1}, \bar{a}_{k+1}) - \sum_{i=1}^k (\bar{a}_{k+1}, \bar{b}_i)^2]^{\frac{1}{2}}} \quad (i=1, 2, \dots, k). \quad (2.5.15)$$

С этой целью введем в рассмотрение угол $\widehat{(\bar{a}_{k+1}, P_k)}$ между вектором \bar{a}_{k+1} и пространством P_k , натянутым на векторы $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_k$. Косинус этого угла определим по формуле

$$\cos(\bar{a}_{k+1}, P_k) = \max_{\bar{z} \in P_k} \cos(\bar{a}_{k+1}, \bar{z}).$$

Так как $\bar{z} \in P_k$, а $\bar{b}_1, \bar{b}_2, \dots, \bar{b}_k$ — базис пространства P_k , то для вектора \bar{z} верно представление

$$\bar{z} = \sum_{i=1}^k \alpha_i \bar{b}_i,$$

где α_i — некоторые числа. Таким образом,

$$\begin{aligned} \cos(\bar{a}_{k+1}, P_k) &= \max_{\bar{z} \in P_k} \cos(\bar{a}_{k+1}, \bar{z}) = \max_{\bar{z} \in P_k} \frac{(\bar{a}_{k+1}, \bar{z})}{\|\bar{a}_{k+1}\|_{III} \cdot \|\bar{z}\|_{III}} = \\ &= \max_{\alpha_i} \frac{(\bar{a}_{k+1}, \sum_{i=1}^k \alpha_i \bar{b}_i)}{\sqrt{(\bar{a}_{k+1}, \bar{a}_{k+1}) \left[\left(\sum_{i=1}^k \alpha_i \bar{b}_i, \sum_{i=1}^k \alpha_i \bar{b}_i \right) \right]^{\frac{1}{2}}}} = \\ &= \frac{1}{\sqrt{(\bar{a}_{k+1}, \bar{a}_{k+1})}} \max_{\alpha_i} \frac{\sum_{i=1}^k \alpha_i (\bar{a}_{k+1}, \bar{b}_i)}{\sqrt{\sum_{i=1}^n \alpha_i^2}}. \end{aligned} \quad (2.5.16)$$

Обозначим $\beta_i = (\bar{a}_{k+1}, \bar{b}_i)$ и введем в рассмотрение вспомогательные векторы $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_k)$; $\bar{\beta} = (\beta_1, \beta_2, \dots, \beta_k)$. Тогда формулу (2.5.16)

можно записать в таком виде:

$$\begin{aligned}
 \cos(\bar{a}_{k+1}, \hat{P}_k) &= \frac{1}{\sqrt{(\bar{a}_{k+1}, \bar{a}_{k+1})}} \max_{\alpha_i} \frac{(\bar{\alpha}, \bar{\beta})}{\sqrt{(\bar{\alpha}, \bar{\alpha})}} = \\
 &= \frac{1}{\sqrt{(\bar{a}_{k+1}, \bar{a}_{k+1})}} \max_{\alpha_i} \frac{(\bar{\alpha}, \bar{\beta}) \cdot \sqrt{(\bar{\beta}, \bar{\beta})}}{\sqrt{(\bar{\alpha}, \bar{\alpha})} \sqrt{(\bar{\beta}, \bar{\beta})}} = \\
 &= \frac{\sqrt{(\bar{\beta}, \bar{\beta})}}{\sqrt{(\bar{a}_{k+1}, \bar{a}_{k+1})}} \max_{\alpha_i} \cos(\bar{\alpha}, \bar{\beta}) = \frac{\sqrt{(\bar{\beta}, \bar{\beta})}}{\sqrt{(\bar{a}_{k+1}, \bar{a}_{k+1})}}, \quad (2.5.17)
 \end{aligned}$$

ибо

$$\max_{\alpha_i} \cos(\bar{\alpha}, \bar{\beta}) = 1.$$

Используя формулы (2.5.15) и (2.5.17) для $\text{ctg}^2(\bar{a}_{k+1}, \hat{P}_k)$, получим такое выражение:

$$\text{ctg}^2(\bar{a}_{k+1}, \hat{P}_k) = \sum_{i=1}^h \gamma_i^2. \quad (2.5.18)$$

Отсюда видно, что если среди чисел $\text{ctg}^2(\bar{a}_{k+1}, \hat{P}_k)$ ($k=1, 2, \dots, n$) есть большие числа, то большими по модулю будут и некоторые из чисел γ_i . В этом случае величины ошибок ε_{k+1j} могут стать значительными по сравнению с ошибками ε_{ij} , полученными на предыдущих шагах. Это в свою очередь означает, что может сильно возрасти погрешность \bar{e}_{k+1} .

Заметим, что числа $\text{ctg}^2(\bar{a}_{k+1}, \hat{P}_k)$ будут большими в том случае, когда угол $(\bar{a}_{k+1}, \hat{P}_k)$ близок к нулю. Близость же указанного угла к нулю означает, что векторы $\bar{a}_1, \bar{a}_2, \dots, \bar{a}_{k+1}$ почти линейно зависимы.

Итак, метод ортогонализации может оказаться неустойчивым к ошибкам округления при решении систем уравнений с матрицами, имеющими две или более почти линейно зависимых строк.

муле (2.5.20) итераций. При $\bar{u}_{k+1} = \bar{u}_{k+1}^{(1)}$ получим обычную схему метода ортогонализации. Как мы видели, в этом случае метод может оказаться неустойчивым. Значительно лучшие результаты получаются при $\bar{u}_{k+1} = \bar{u}_{k+1}^{(2)}$. Вообще же число итераций обуславливается требованиями точности вычисления вектора \bar{u}_{k+1} .

Отметим, наконец, что алгоритм Уилкинсона направлен на более точное вычисление вектора \bar{u}_{k+1} . После того как этот вектор вычислен, дальнейшие вычисления проводятся по схеме метода ортогонализации, т. е. вычисляется вектор \bar{b}_{k+1} и т. д.

2.5.3. Метод сопряженных градиентов

Этот метод предназначен для решения системы линейных алгебраических уравнений

$$A\bar{x} = \bar{f} \quad (2.5.25)$$

с вещественной симметрической положительно определенной матрицей.

В методе сопряженных градиентов отыскание решения системы (2.5.25) связывается с задачей минимизации следующего функционала:

$$F(\bar{x}) = (A\bar{x}, \bar{x}) - 2(\bar{f}, \bar{x}), \quad (2.5.26)$$

являющегося квадратичной функцией относительно x_1, x_2, \dots, x_n . Дело в том, что решение системы (2.5.25) — вектор $\bar{x}^{(*)} = A^{-1}\bar{f}$ — доставляет минимум функционалу (2.5.26) на множестве векторов из вещественного векторного пространства. Действительно, из (2.5.25), (2.5.26) следует

$$\begin{aligned} F(\bar{x}) - F(\bar{x}^{(*)}) &= (A\bar{x}, \bar{x}) - 2(\bar{f}, \bar{x}) - (A\bar{x}^{(*)}, \bar{x}^{(*)}) + 2(\bar{f}, \bar{x}^{(*)}) = (A\bar{x}, \bar{x}) - 2(A\bar{x}^{(*)}, \bar{x}) - \\ &- (A\bar{x}^{(*)}, \bar{x}^{(*)}) + 2(A\bar{x}^{(*)}, \bar{x}^{(*)}) = (A(\bar{x} - \bar{x}^{(*)}), (\bar{x} - \bar{x}^{(*)})) \geq 0, \end{aligned} \quad (2.5.27)$$

так как матрица A положительно определенная. При этом знак равенства в (2.5.27) возможен лишь при $\bar{x} - \bar{x}^{(*)} = 0$, т. е. при $\bar{x} = \bar{x}^{(*)}$. Таким образом, задача нахождения решения системы (2.5.25) сводится к задаче отыскания вектора \bar{x} , доставляющего минимум функционалу $F(\bar{x})$.

Прежде чем переходить к изложению правила для отыскания такого вектора, остановимся кратко на понятии градиента функционала. Пусть $F(\bar{x})$ — некоторый функционал и пусть \bar{y} — произвольный вектор единичной длины с координатами y_1, y_2, \dots, y_n . Производной от функционала F в точке \bar{x} по направлению \bar{y} называется выражение

$$\frac{\partial F(\bar{x})}{\partial \bar{y}} = \lim_{t \rightarrow 0} \frac{F(\bar{x} + t\bar{y}) - F(\bar{x})}{t} = \frac{d}{dt} F(\bar{x} + t\bar{y}) \Big|_{t=0}.$$

Производная $\frac{\partial F(\bar{x})}{\partial \bar{y}}$ характеризует скорость изменения функционала F при изменении «аргумента» в направлении вектора \bar{y} . Имеем далее

$$F(\bar{x} + t\bar{y}) = F(x_1 + ty_1, \dots, x_n + ty_n),$$

поэтому

$$\frac{\partial F(\bar{x})}{\partial \bar{y}} = \frac{d}{dt} F(x_1 + ty_1, \dots, x_n + ty_n) \Big|_{t=0} = \frac{\partial F(\bar{x})}{\partial x_1} y_1 + \dots + \frac{\partial F(\bar{x})}{\partial x_n} y_n = (\bar{z}, \bar{y}),$$

где

$$\bar{z} = (z_1, z_2, \dots, z_n)', \quad z_i = \frac{\partial F(\bar{x})}{\partial x_i}.$$

Вектор \bar{z} называется градиентом функционала $F(\bar{x})$. Из последнего равенства вытекает, что

$$\frac{\partial F(\bar{x})}{\partial \bar{y}} = |\bar{z}| \cos \angle(\bar{z}, \bar{y}),$$

ибо $|\bar{y}| = 1$. Отсюда следует, что

$$-|\bar{z}| \leq \frac{\partial F(\bar{x})}{\partial \bar{y}} \leq |\bar{z}|,$$

причем $\frac{\partial F(\bar{x})}{\partial \bar{y}} = |\bar{z}|$, если направление \bar{y} совпадает с направлением градиента, и $\frac{\partial F(\bar{x})}{\partial \bar{y}} = -|\bar{z}|$, если направление противоположно направлению градиента. Поэтому

направление градиента есть направление наибольшей скорости роста функционала F в данной точке, а направление, противоположное градиенту, есть направление наибольшей скорости убывания. Это последнее направление существенно используется в методе сопряженных градиентов для отыскания минимума функционала $F(\bar{x})$.

Перейдем к изложению схемы метода. Вектор, доставляющий минимум функционалу $F(\bar{x})$, будем находить итерационным способом. Пусть $\bar{x}^{(0)}$ — произвольный начальный вектор. Рассмотрим функционал (2.5.26) и вычислим его градиент. Имеем:

$$\begin{aligned} \frac{\partial F(\bar{x})}{\partial \bar{y}} &= \frac{d}{dt} F(\bar{x} + t\bar{y}) \Big|_{t=0} = \frac{d}{dt} (A(\bar{x} + t\bar{y}) - 2\bar{f}, \bar{x} + t\bar{y}) \Big|_{t=0} = \\ &= \frac{d}{dt} [t^2(A\bar{y}, \bar{y}) - 2t(\bar{f} - A\bar{x}, \bar{y}) + F(\bar{x})]_{t=0} = -2(\bar{f} - A\bar{x}, \bar{y}) = 2(A\bar{x} - \bar{f}, \bar{y}). \end{aligned}$$

Следовательно, градиент $F(\bar{x})$ равен $2A\bar{x} - 2\bar{f}$. Так как в дальнейшем нам важно лишь направление градиента, мы отбрасываем положительный множитель 2 и будем рассматривать вместо градиента функционала $F(\bar{x})$ вектор $A\bar{x} - \bar{f}$. Вектор, имеющий в точке $\bar{x}^{(0)}$ направление, противоположное градиенту, обозначим через $\bar{r}^{(0)}$, т. е.

$$\bar{r}^{(0)} = \bar{f} - A\bar{x}^{(0)}. \quad (2.5.28)$$

Заметим, что в направлении этого вектора, который мы будем называть также вектором невязок системы, скорость убывания функционала $F(\bar{x})$ в точке $\bar{x}^{(0)}$ наибольшая. Будем теперь двигаться из точки $\bar{x}^{(0)}$ в направлении вектора $\bar{r}^{(0)}$ до тех пор, пока функция $F(\bar{x}^{(0)} + \alpha \bar{r}^{(0)})$ достигнет своего минимального значения в этом направлении. Это будет при $\frac{d}{d\alpha} F(\bar{x}^{(0)} + \alpha \bar{r}^{(0)}) = 0$, т. е. при

$$\alpha_0 = \frac{(\bar{r}^{(0)}, \bar{r}^{(0)})}{(\bar{r}^{(0)}, A\bar{r}^{(0)})}. \quad (2.5.29)$$

Здесь $(\bar{r}^{(0)}, A\bar{r}^{(0)}) > 0$ при любых $\bar{r}^{(0)} \neq 0$ в силу предположений относительно матрицы A . Если $\bar{r}^{(0)} = 0$, то из (2.5.28) видно, что $\bar{x}^{(0)}$ совпадает с решением и никаких дальнейших вычислений проводить не следует. За новое приближение к решению при $\bar{r}^{(0)} \neq 0$ принимаем вектор

$$\bar{x}^{(1)} = \bar{x}^{(0)} + \alpha_0 \bar{r}^{(0)}. \quad (2.5.30)$$

Отметим еще, что вектор невязок $\bar{r}^{(0)}$ имеет направление нормали к поверхности $F(\bar{x}) = F(\bar{x}^{(0)})$ в точке $\bar{x} = \bar{x}^{(0)}$, ибо направление быстрого изменения функции $F(\bar{x})$ в этой точке совпадает с направлением нормали.

Следующее приближение $\bar{x}^{(2)}$ находится так. Обозначим через Γ_k гиперплоскость k измерений (способы задания Γ_k для каждого k будем указывать ниже) и проведем через точку $\bar{x}^{(1)}$ гиперплоскость Γ_{n-1}

$$(A\bar{r}^{(0)}, \bar{x} - \bar{x}^{(1)}) = 0.$$

Обозначим через $\bar{r}^{(1)}$ новую невязку системы

$$\bar{r}^{(1)} = \bar{f} - A\bar{x}^{(1)} = \bar{r}^{(0)} - \alpha_0 A\bar{r}^{(0)}. \quad (2.5.31)$$

Вектор $\bar{r}^{(1)}$ направлен по нормали к поверхности $F(\bar{x}) = F(\bar{x}^{(1)})$ в точке $\bar{x} = \bar{x}^{(1)}$, а вектор $\bar{r}^{(0)}$ параллелен касательной плоскости в этой точке. Значит, $\bar{r}^{(0)}$ и $\bar{r}^{(1)}$ ортогональны, т. е.

$$(\bar{r}^{(0)}, \bar{r}^{(1)}) = 0. \quad (2.5.32)$$

Ортогональность $\bar{r}^{(0)}$ и $\bar{r}^{(1)}$ следует также из формул (2.5.28) — (2.5.31), в чем легко убедиться. Гиперплоскость Γ_{n-1} проходит через точку $\bar{x}^{(*)} = A^{-1}\bar{f}$, ибо

$$(A\bar{r}^{(0)}, A^{-1}\bar{f} - \bar{x}^{(1)}) = (\bar{r}^{(0)}, \bar{f} - A\bar{x}^{(1)}) = (\bar{r}^{(0)}, \bar{r}^{(1)}) = 0.$$

Теперь нам известно, что решение системы лежит в гиперплоскости Γ_{n-1} , проходящей через точку $\bar{x}^{(1)}$. Однако нам не известно направление, двигаясь по которому в гиперплоскости Γ_{n-1} можно достичь точки $\bar{x}^{(*)}$. Пока у нас нет достаточных сведений для определения этого направления и все, что мы можем сделать, это определить некоторый вектор $\bar{p}^{(1)}$, лежащий в Γ_{n-1} , и затем двигаться из точки $\bar{x}^{(1)}$ в направлении этого вектора до тех пор, пока функция $F(\bar{x}^{(1)} + \alpha \bar{p}^{(1)})$ не достигнет минимума. По построению вектор

$\bar{r}^{(1)} + \beta \bar{r}^{(0)}$ параллелен некоторой нормальной плоскости к поверхности $F(\bar{x}) = F(\bar{x}^{(1)})$ в точке $\bar{x}^{(1)}$ при любом β . Выберем β таким образом, чтобы этот вектор лежал в Γ_{n-1} , т. е. был ортогонален к $A\bar{r}^{(0)}$. Это дает

$$(\bar{r}^{(1)} + \beta_0 \bar{r}^{(0)}, A\bar{r}^{(0)}) = (\bar{r}^{(1)}, A\bar{r}^{(0)}) + \beta_0 (\bar{r}^{(0)}, A\bar{r}^{(0)}) = 0.$$

Отсюда

$$\beta_0 = - \frac{(\bar{r}^{(1)}, A\bar{r}^{(0)})}{(\bar{r}^{(0)}, A\bar{r}^{(0)})}. \quad (2.5.33)$$

Таким образом, по смыслу задачи в качестве вектора $\bar{p}^{(1)}$ можно принять вектор $\bar{r}^{(1)} + \beta_0 \bar{r}^{(0)}$, лежащий в Γ_{n-1} :

$$\bar{p}^{(1)} = \bar{r}^{(1)} + \beta_0 \bar{r}^{(0)}. \quad (2.5.34)$$

Отметим, что этот вектор имеет направление нормали к сечению поверхности $F(\bar{x}) = F(\bar{x}^{(1)})$ гиперплоскостью Γ_{n-1} в точке $\bar{x}^{(1)}$. Далее, из условия $\frac{d}{d\alpha} F(\bar{x}^{(1)} + \alpha \bar{p}^{(1)}) = 0$ получим

$$\alpha_1 = \frac{(\bar{r}^{(1)}, \bar{p}^{(1)})}{(\bar{p}^{(1)}, A\bar{p}^{(1)})}. \quad (2.5.35)$$

В качестве второго приближения к решению системы примем вектор $\bar{x}^{(2)}$:

$$\bar{x}^{(2)} = \bar{x}^{(1)} + \alpha_1 \bar{p}^{(1)}. \quad (2.5.36)$$

Укажем еще правило для вычисления вектора $\bar{x}^{(3)}$, после чего станет ясной вся схема метода сопряженных градиентов.

Как и векторы $\bar{r}^{(0)}$, $\bar{r}^{(1)}$, вектор невязок

$$\bar{r}^{(2)} = \bar{f} - A\bar{x}^{(2)} = \bar{r}^{(1)} - \alpha_1 A\bar{p}^{(1)} \quad (2.5.37)$$

имеет направление нормали к соответствующей поверхности $F(\bar{x}) = F(\bar{x}^{(2)})$ в точке $\bar{x}^{(2)}$. Покажем, что $\bar{r}^{(2)}$ ортогонален к $\bar{r}^{(0)}$ и $\bar{r}^{(1)}$. Действительно, используя (2.5.32) — (2.5.37), получим:

$$(\bar{r}^{(2)}, \bar{r}^{(0)}) = (\bar{r}^{(1)} - \alpha_1 A\bar{p}^{(1)}, \bar{r}^{(0)}) = -\alpha_1 (A\bar{p}^{(1)}, \bar{r}^{(0)}) = -\alpha_1 (\bar{p}^{(1)}, A\bar{r}^{(0)}) = 0;$$

$$(\bar{r}^{(2)}, \bar{r}^{(1)}) = (\bar{r}^{(1)} - \alpha_1 A\bar{p}^{(1)}, \bar{p}^{(1)} - \beta_0 \bar{r}^{(0)}) = (\bar{r}^{(1)}, \bar{p}^{(1)}) - \alpha_1 (\bar{p}^{(1)}, A\bar{p}^{(1)}) = 0.$$

Рассмотрим гиперплоскость Γ_{n-2}

$$(A\bar{r}^{(0)}, \bar{x} - \bar{x}^{(1)}) = 0, \quad (A\bar{p}^{(1)}, \bar{x} - \bar{x}^{(2)}) = 0,$$

проходящую через точку $\bar{x}^{(2)}$. На Γ_{n-2} лежит и точка $\bar{x}^{(*)}$, ибо $(A\bar{r}^{(0)}, \bar{x}^{(*)} - \bar{x}^{(1)}) = 0$, так как $\bar{x}^{(*)} \in \Gamma_{n-1}$, а

$$(A\bar{p}^{(1)}, \bar{x}^{(*)} - \bar{x}^{(2)}) = (\bar{p}^{(1)}, A'(\bar{x}^{(*)} - \bar{x}^{(2)})) = (\bar{p}^{(1)}, A'A^{-1}\bar{f} - A'\bar{x}^{(2)}) = (\bar{r}^{(1)} + \beta_0 \bar{r}^{(0)}, \bar{r}^{(2)}) = 0,$$

Теперь мы находимся в тех же условиях, что и при нахождении $\bar{x}^{(2)}$, а именно: нам известно приближение $\bar{x}^{(2)}$, гиперплоскость Γ_{n-2} , проходящая через $\bar{x}^{(2)}$ и через искомое решение $\bar{x}^{(*)}$. Далее будем действовать аналогично предыдущему. Вектор $\bar{r}^{(2)} + \beta \bar{p}^{(1)}$ параллелен Γ_{n-1} при любом β , ибо

$$(\bar{r}^{(2)} + \beta \bar{p}^{(1)}, A\bar{r}^{(0)}) = (\bar{r}^{(2)}, A\bar{r}^{(0)}) + \beta (\bar{p}^{(1)}, A\bar{r}^{(0)}) = \left(\bar{r}^{(2)}, \frac{1}{\alpha_0} (\bar{r}^{(0)} - \bar{r}^{(1)}) \right) = 0.$$

Выберем β так, чтобы $\bar{r}^{(2)} + \beta \bar{p}^{(1)}$ был параллелен Γ_{n-2} , т. е. потребуем ортогональности этого вектора к вектору $A\bar{p}^{(1)}$. Это дает следующее условие для определения β_1 :

$$(\bar{r}^{(2)} + \beta_1 \bar{p}^{(1)}, A\bar{p}^{(1)}) = (\bar{r}^{(2)}, A\bar{p}^{(1)}) + \beta_1 (\bar{p}^{(1)}, A\bar{p}^{(1)}) = 0$$

или

$$\beta_1 = - \frac{(\bar{r}^{(2)}, A\bar{p}^{(1)})}{(\bar{p}^{(1)}, A\bar{p}^{(1)})}. \quad (2.5.38)$$

Вектор

$$\bar{p}^{(2)} = \bar{r}^{(2)} + \beta_1 \bar{p}^{(1)} \quad (2.5.39)$$

будет иметь направление нормали к сечению поверхности $F(\bar{x}) = F(\bar{x}^{(2)})$ гиперплоскостью Γ_{n-2} в точке $\bar{x}^{(2)}$. Минимизируя функцию $F(\bar{x}^{(2)} + \alpha \bar{p}^{(2)})$, получим для α_2 такое выражение:

$$\alpha_2 = \frac{(\bar{r}^{(2)}, \bar{p}^{(2)})}{(\bar{p}^{(2)}, A\bar{p}^{(2)})}. \quad (2.5.40)$$

Значит, вектор

$$\bar{x}^{(3)} = \bar{x}^{(2)} + \alpha_2 \bar{p}^{(2)} \quad (2.5.41)$$

будет новым приближением к \bar{x}^* . По этому вектору вычисляем новый вектор невязок

$$\bar{r}^{(3)} = \bar{f} - A\bar{x}^{(3)} = \bar{r}^{(2)} - \alpha_2 A\bar{p}^{(2)}$$

и продолжаем процесс по аналогии с предыдущим. В результате мы получим последовательности векторов $\{\bar{x}^{(k)}\}$, $\{\bar{r}^{(k)}\}$, $\{\bar{p}^{(k)}\}$ и чисел $\{\alpha_k\}$, $\{\beta_k\}$, определяемые следующими рекуррентными соотношениями:

$$\left. \begin{aligned} \bar{p}^{(0)} &= \bar{r}^{(0)} = \bar{f} - A\bar{x}^{(0)}, & \alpha_k &= \frac{(\bar{r}^{(k)}, \bar{p}^{(k)})}{(\bar{p}^{(k)}, A\bar{p}^{(k)})}, \\ \bar{x}^{(k+1)} &= \bar{x}^{(k)} + \alpha_k \bar{p}^{(k)}, & \bar{r}^{(k+1)} &= \bar{f} - A\bar{x}^{(k+1)} = \bar{r}^{(k)} - \alpha_k A\bar{p}^{(k)}, \\ \beta_k &= - \frac{(\bar{r}^{(k+1)}, A\bar{p}^{(k)})}{(\bar{p}^{(k)}, A\bar{p}^{(k)})}, & \bar{p}^{(k+1)} &= \bar{r}^{(k+1)} + \beta_k \bar{p}^{(k)}, \\ k &= 0, 1, \dots, s, & s &\leq n. \end{aligned} \right\} \quad (2.5.42)$$

Докажем конечность определяемого формулами (2.5.42) алгоритма, для чего установим факт, что $\bar{x}^{(k)} = \bar{x}^{(*)}$ при некотором $k \leq n$.

Покажем сначала, что

$$(\bar{r}^{(i)}, \bar{p}^{(j)}) = 0 \quad \text{при } i > j, \quad (2.5.43)$$

$$(\bar{r}^{(i)}, \bar{r}^{(j)}) = 0 \quad \text{при } i \neq j. \quad (2.5.44)$$

По построению

$$(\bar{p}^{(i)}, A\bar{p}^{(j)}) = (A^* \bar{p}^{(i)}, \bar{p}^{(j)}) = (A\bar{p}^{(i)}, \bar{p}^{(j)}) = 0$$

при $i \neq j$. Далее при $i > j$

$$(\bar{r}^{(i)}, \bar{p}^{(j)}) = (\bar{r}^{(i-1)} - \alpha_{i-1} A\bar{p}^{(i-1)}, \bar{p}^{(j)}) = (\bar{r}^{(i-1)}, \bar{p}^{(j)}) - \alpha_{i-1} (A\bar{p}^{(i-1)}, \bar{p}^{(j)}).$$

Здесь правая часть равна нулю, если $i = j+1$, в силу определения α_{i-1} . Если $i > j+1$, то $(A\bar{p}^{(i-1)}, \bar{p}^{(j)}) = 0$ и, значит,

$$(\bar{r}^{(i)}, \bar{p}^{(j)}) = (\bar{r}^{(i-1)}, \bar{p}^{(j)}). \quad (2.5.45)$$

Отсюда получим

$$\begin{aligned} (\bar{r}^{(i)}, \bar{p}^{(j)}) &= (\bar{r}^{(i-1)}, \bar{p}^{(j)}) = (\bar{r}^{(i-2)}, \bar{p}^{(j)}) = \dots = (\bar{r}^{(j+1)}, \bar{p}^{(j)}) = \\ &= (\bar{r}^{(j)} - \alpha_j A\bar{p}^{(j)}, \bar{p}^{(j)}) = (\bar{r}^{(j)}, \bar{p}^{(j)}) - \alpha_j (A\bar{p}^{(j)}, \bar{p}^{(j)}) = 0 \end{aligned}$$

в силу определения α_j . Итак, соотношения (2.5.43) доказаны. Для доказательства (2.5.44) предположим, например, что $i > j$. Тогда

$$(\bar{r}^{(i)}, \bar{r}^{(j)}) = (\bar{r}^{(i)}, \bar{p}^{(j)} - \beta_{j-1} \bar{p}^{(j-1)}) = (\bar{r}^{(i)}, \bar{p}^{(j)}) - \beta_{j-1} (\bar{r}^{(i)}, \bar{p}^{(j-1)}) = 0$$

в силу (2.5.43). Понятно, что верно и $(\bar{r}^{(i)}, \bar{r}^{(j)}) = 0$ при $i < j$.

Таким образом, система векторов $\bar{r}^{(0)}, \bar{r}^{(1)}, \dots, \bar{r}^{(k)}$ ортогональна. Но поскольку в n -мерном векторном пространстве не может быть более n взаимно ортогональных векторов, то на некотором шаге $k \leq n$ получим $\bar{r}^{(k)} = 0$. Значит, $\bar{f} - A\bar{x}^{(k)} = 0$ и вектор $\bar{x}^{(k)}$ будет искомым решением системы уравнений $A\bar{x} = \bar{f}$.

Вся схема метода сопряженных градиентов определяется формулами (2.5.42). Реализуются эти формулы просто. В процессе вычислений контроль точности вычисления вектора $\bar{x}^{(k)}$ можно проводить путем оценки вектора невязок $\bar{r}^{(k)}$ в какой-либо метрике. Порядок систем, решаемых на ЭВМ, в основном зависит от объема числовой информации, необходимой для определения элементов матрицы A . Это объясняется тем, что основной операцией в методе является многократное вычисление произведений матрицы A на векторы $\bar{p}^{(k)}$. Поэтому метод сопряженных градиентов целесообразно использовать для решения систем уравнений, в которых матрица A имеет много нулевых элементов. В этом случае умножение A на $\bar{p}^{(k)}$ для ЭВМ можно организовать так, чтобы в арифметических операциях участвовали только ненулевые элементы матрицы.

Методу сопряженных градиентов свойствен и некоторый недостаток. Проводимый в этом методе процесс ортогонализации может оказаться неустойчивым к ошибкам округления, как это было в случае метода ортогонализации. Чтобы ослабить неустойчивость, надо время от времени, вычисляя вектор $\bar{r}^{(k)}$ по формуле $\bar{r}^{(k)} = \bar{r}^{(k-1)} - \alpha_k A\bar{p}^{(k-1)}$, проводить затем еще вычисления и по формуле $\bar{r}^{(k)} = \bar{f} - A\bar{x}^{(k)}$, и при расхождении брать второй результат.

Отметим, что метод сопряженных градиентов может быть распространен и на системы уравнений с произвольной невырожденной матрицей.

2.5.4. Вариант метода сопряженных градиентов

Рассматриваемый ниже вариант метода требует несколько больших вычислений, чем это было в предыдущем случае. Зато он менее чувствителен к ошибкам округления, что часто позволяет получить хорошее приближение к решению за меньшее число шагов по сравнению с методом сопряженных градиентов.

Суть варианта в следующем. Как и ранее, будем разыскивать минимум функционала $F(\bar{x})$ на некотором множестве X векторов \bar{x} , определяемых формулой

$$\bar{x} = \bar{y}^{(i)} + \gamma \bar{\eta}^{(i)} + \delta \bar{\rho}^{(i)}, \quad (2.5.46)$$

где $\bar{y}^{(i)}$, $\bar{\eta}^{(i)}$, $\bar{\rho}^{(i)}$ — некоторые векторы, γ и δ — числа. Укажем правило для их определения.

Вектор $\bar{y}^{(i)}$ будет обозначать i -е последовательное приближение к решению системы $A\bar{x} = \bar{f}$, в качестве $\bar{\eta}^{(i)}$ и $\bar{\rho}^{(i)}$ возьмем следующие векторы:

$$\bar{\eta}^{(i)} = \bar{y}^{(i)} - \bar{y}^{(i-1)}, \quad \bar{\rho}^{(i)} = \bar{f} - A\bar{y}^{(i)}. \quad (2.5.47)$$

Будем считать, что два начальных приближения $\bar{y}^{(0)}$ и $\bar{y}^{(1)}$ выбраны таким образом, что выполняется условие

$$(\bar{\rho}^{(1)}, \bar{\eta}^{(1)}) = 0. \quad (2.5.48)$$

Пусть мы уже вычислили векторы $\bar{y}^{(k)}$, $\bar{\eta}^{(k)}$, $\bar{\rho}^{(k)}$ ($k=0, 1, \dots, i$) такие, что справедливы равенства

$$(\bar{\rho}^{(k)}, \bar{\eta}^{(k)}) = 0 \quad (k=1, \dots, i). \quad (2.5.49)$$

Чтобы вычислить минимум функционала $F(\bar{x})$ на множестве X , приравняем нулю частные производные от $F(\bar{y}^{(i)} + \gamma \bar{\eta}^{(i)} + \delta \bar{\rho}^{(i)})$ по γ и по δ . Это даст следующую систему уравнений для определения γ и δ :

$$\left. \begin{aligned} (A(\bar{y}^{(i)} + \gamma \bar{\eta}^{(i)} + \delta \bar{\rho}^{(i)}) - \bar{f}, \bar{\eta}^{(i)}) &= 0, \\ (A(\bar{y}^{(i)} + \gamma \bar{\eta}^{(i)} + \delta \bar{\rho}^{(i)}) - \bar{f}, \bar{\rho}^{(i)}) &= 0 \end{aligned} \right\}$$

или, с учетом (2.5.49) при $k=i$,

$$\left. \begin{aligned} \gamma(\bar{\eta}^{(i)}, A\bar{\eta}^{(i)}) + \delta(\bar{\rho}^{(i)}, A\bar{\eta}^{(i)}) &= 0, \\ \gamma(\bar{\rho}^{(i)}, A\bar{\eta}^{(i)}) + \delta(\bar{\rho}^{(i)}, A\bar{\rho}^{(i)}) &= (\bar{\rho}^{(i)}, \bar{\rho}^{(i)}). \end{aligned} \right\} \quad (2.5.50)$$

Обозначим через γ_i , δ_i решение этой системы:

$$\left. \begin{aligned} \gamma_i &= \frac{-(\bar{\rho}^{(i)}, A\bar{\eta}^{(i)}) (\bar{\rho}^{(i)}, \bar{\rho}^{(i)})}{(\bar{\eta}^{(i)}, A\bar{\eta}^{(i)}) (\bar{\rho}^{(i)}, A\bar{\rho}^{(i)}) - (\bar{\rho}^{(i)}, A\bar{\eta}^{(i)})^2}, \\ \delta_i &= \frac{(\bar{\eta}^{(i)}, A\bar{\eta}^{(i)}) (\bar{\rho}^{(i)}, \bar{\rho}^{(i)})}{(\bar{\eta}^{(i)}, A\bar{\eta}^{(i)}) (\bar{\rho}^{(i)}, A\bar{\rho}^{(i)}) - (\bar{\rho}^{(i)}, A\bar{\eta}^{(i)})^2} \end{aligned} \right\} \quad (2.5.51)$$

Отметим, что определитель системы (2.5.50) отличен от нуля, ибо в силу неравенства Коши — Буняковского

$$(\bar{\rho}^{(i)}, A\bar{\eta}^{(i)})^2 \leq (\bar{\eta}^{(i)}, A\bar{\eta}^{(i)}) (\bar{\rho}^{(i)}, A\bar{\rho}^{(i)}),$$

причем равенство здесь имеет место тогда и только тогда, когда векторы $\bar{\rho}^{(i)}$ и $\bar{\eta}^{(i)}$ коллинеарны. За $(i+1)$ -е приближение к решению системы примем вектор

$$\bar{y}^{(i+1)} = \bar{y}^{(i)} + \gamma_i \bar{\eta}^{(i)} + \delta_i \bar{\rho}^{(i)}. \quad (2.5.52)$$

Покажем, что равенство (2.5.49) выполняется и при $k=i+1$. Действительно, в силу (2.5.47), (2.5.50) и (2.5.52)

$$(\bar{\rho}^{(i+1)}, \bar{\eta}^{(i)}) = 0, \quad (\bar{\rho}^{(i+1)}, \bar{\rho}^{(i)}) = 0. \quad (2.5.53)$$

Значит,

$$(\bar{\rho}^{(i+1)}, \bar{\eta}^{(i+1)}) = (\bar{\rho}^{(i+1)}, \bar{y}^{(i+1)} - \bar{y}^{(i)}) = (\bar{\rho}^{(i+1)}, \gamma_i \bar{\eta}^{(i)} + \delta_i \bar{\rho}^{(i)}) = 0. \quad (2.5.54)$$

Докажем теперь теорему о сходимости и конечности варианта метода сопряженных градиентов.

Теорема 1. Если начальные приближения $\bar{y}^{(0)}$ и $\bar{y}^{(1)}$ к решению $\bar{x}^{(*)}$ системы $A\bar{x} = \bar{f}$ выбраны таким образом, что

$$\bar{y}^{(0)} = \bar{x}^{(0)}, \quad \bar{y}^{(1)} = \bar{x}^{(1)}, \quad (2.5.55)$$

то при всех i будет иметь место равенство

$$\bar{y}^{(i)} = \bar{x}^{(i)}, \quad (2.5.56)$$

где $\bar{x}^{(i)}$ — вектор, определяемый формулами (2.5.42) и обладающий свойством $\bar{x}^{(k)} = \bar{x}^{(*)}$ при некотором $k \leq n$ (n — порядок системы).

Доказательство. Сначала проверим, что условие (2.5.49) выполняется для $k=1$, если векторы $\bar{y}^{(0)}$ и $\bar{y}^{(1)}$ определяются по формуле (2.5.55). Действительно,

$$\bar{\rho}^{(1)} = \bar{f} - A\bar{x}^{(1)} = \bar{r}^{(1)}, \quad \bar{\eta}^{(1)} = \bar{y}^{(1)} - \bar{y}^{(0)} = \bar{x}^{(1)} - \bar{x}^{(0)} = \alpha_0 \bar{r}^{(0)}$$

и, значит,

$$(\bar{\rho}^{(1)}, \bar{\eta}^{(1)}) = \alpha_0 (\bar{r}^{(1)}, \bar{r}^{(0)}) = 0$$

в силу (2.5.32). Равенство $\bar{y}^{(k)} = \bar{x}^{(k)}$ будем доказывать по индукции. Предположим, что это равенство имеет место при $k=0, 1, \dots, i$, и докажем его справедливость при $k=i+1$. Имеем

$$\bar{\rho}^{(i)} = \bar{f} - A\bar{y}^{(i)} = \bar{f} - A\bar{x}^{(i)} = \bar{r}^{(i)},$$

ибо по предположению $\bar{y}^{(i)} = \bar{x}^{(i)}$. Значит,

$$\bar{y}^{(i+1)} = \bar{x}^{(i)} + \gamma_i \bar{\eta}^{(i)} + \delta_i \bar{r}^{(i)}, \quad \bar{x}^{(i+1)} = \bar{x}^{(i)} + \alpha_i \bar{r}^{(i)}.$$

Для вектора $\bar{\rho}^{(i)}$, используя (2.5.42), получим

$$\bar{\rho}^{(i)} = \bar{r}^{(i)} + \beta_{i-1} \bar{\rho}^{(i-1)} = \bar{r}^{(i)} + \frac{\beta_{i-1}}{\alpha_{i-1}} (\bar{x}^{(i)} - \bar{x}^{(i-1)}) = \bar{r}^{(i)} + \frac{\beta_{i-1}}{\alpha_{i-1}} \bar{\eta}^{(i)},$$

Поэтому

$$\bar{x}^{(i+1)} = \bar{x}^{(i)} + \alpha_i \bar{r}^{(i)} + \frac{\alpha_i \beta_{i-1}}{\alpha_{i-1}} \bar{\eta}^{(i)}.$$

Докажем теперь коллинеарность векторов

$$\gamma_i \bar{\eta}^{(i)} + \delta_i \bar{r}^{(i)} \quad \text{и} \quad \frac{\alpha_i \beta_{i-1}}{\alpha_{i-1}} \bar{\eta}^{(i)} + \alpha_i \bar{r}^{(i)}. \quad (2.5.57)$$

Из (2.5.42) и (2.5.51) имеем

$$\begin{aligned} \frac{\gamma_i}{\delta_i} &= - \frac{(\bar{p}^{(i)}, A\bar{\eta}^{(i)})}{(\bar{\eta}^{(i)}, A\bar{\eta}^{(i)})} = - \frac{(\bar{r}^{(i)}, A\alpha_{i-1}\bar{p}^{(i-1)})}{(\alpha_{i-1}\bar{p}^{(i-1)}, A\alpha_{i-1}\bar{p}^{(i-1)})} = \\ &= \frac{\alpha_{i-1}}{\alpha_{i-1}^2} \cdot \frac{(\bar{r}^{(i)}, A\bar{p}^{(i-1)})}{(\bar{p}^{(i-1)}, A\bar{p}^{(i-1)})} = \frac{\alpha_{i-1}}{\alpha_{i-1}^2} \beta_{i-1} = \frac{\beta_{i-1}}{\alpha_{i-1}} = \frac{\alpha_i \beta_{i-1}}{\alpha_{i-1}} \cdot \frac{1}{\alpha_i}, \end{aligned}$$

что и доказывает коллинеарность векторов (2.5.57).

По построению вектор $\bar{y}^{(i+1)}$ дает минимум функционалу F в плоскости, проходящей через точку $\bar{x}^{(i)}$ и натянутой на векторы $\bar{\eta}^{(i)}$ и $\bar{r}^{(i)}$. Но этот же минимум, по ранее доказанному предположению в методе сопряженных градиентов, лежит на прямой, проходящей через точку $\bar{x}^{(i)}$ в направлении вектора $\bar{p}^{(i)}$, и достигается на векторе $\bar{x}^{(i+1)}$. А это и означает, что

$$\bar{y}^{(i+1)} = \bar{x}^{(i+1)}.$$

Таким образом, равенство $\bar{y}^{(k)} = \bar{x}^{(k)}$ имеет место и при $k = i+1$, что и доказывает справедливость равенства (2.5.56) при всех i .

Теорема доказана.

В заключение изложим порядок вычислений в рассматриваемом варианте метода сопряженных градиентов.

При решении системы уравнений $A\bar{x} = \bar{f}$ с симметрической положительно определенной матрицей выбираем сначала вектор $\bar{x}^{(0)}$ — некоторое нулевое приближение к решению $\bar{x}^{(*)}$ — и полагаем $\bar{y}^{(0)} = \bar{x}^{(0)}$. Затем вычисляем вектор $\bar{r}^{(0)} = \bar{f} - A\bar{x}^{(0)}$ и число

$$\alpha_0 = \frac{(\bar{r}^{(0)}, \bar{r}^{(0)})}{(\bar{r}^{(0)}, A\bar{r}^{(0)})},$$

а по ним — вектор $\bar{x}^{(1)} = \bar{x}^{(0)} + \alpha_0 \bar{r}^{(0)}$. Полагаем $\bar{y}^{(1)} = \bar{x}^{(1)}$. Далее, по формулам (2.5.47) при $i=1$ вычисляем $\bar{\eta}^{(1)}$ и $\bar{\rho}^{(1)}$, а по формулам (2.5.51) при $i=1$ — числа γ_1 и δ_1 . Второе приближение к решению — вектор $\bar{y}^{(2)}$ — вычисляем по формуле (2.5.52) при $i=1$. Для вычисления вектора $\bar{y}^{(3)}$ используем формулы (2.5.47), (2.5.51) и (2.5.52) при $i=2$ и т. д. Контроль точности вычисления векторов $\bar{y}^{(i)}$ можно осуществлять путем оценки векторов $\bar{\eta}^{(i)}$ и $\bar{\rho}^{(i)}$ в некоторой метрике, например в евклидовой.

2.5.5. Метод скорейшего спуска

В основе метода, как и в методе сопряженных градиентов, лежит идея нахождения вектора, доставляющего минимум функционалу $F(\bar{x})$ (2.5.26). Как мы видели, такой вектор является решением системы уравнений $A\bar{x}=\bar{f}$ с положительно определенной симметрической матрицей.

Метод имеет следующую вычислительную схему. Исходя из некоторого начального приближения $\bar{x}^{(0)}$ к решению системы $\bar{x}^{(*)}$, вычисляем по такому же правилу, как и в методе сопряженных градиентов, вектор $\bar{r}^{(0)}=\bar{f}-A\bar{x}^{(0)}$, число

$$\alpha_0 = \frac{(\bar{r}^{(0)}, \bar{r}^{(0)})}{(\bar{r}^{(0)}, A\bar{r}^{(0)})}$$

и следующее приближение — вектор

$$\bar{x}^{(1)} = \bar{x}^{(0)} + \alpha_0 \bar{r}^{(0)}.$$

Вектор $\bar{x}^{(2)}$ вычисляем из условия минимума функции $F(\bar{x}^{(1)} + \alpha \bar{r}^{(1)})$, где $\bar{r}^{(1)} = \bar{f} - A\bar{x}^{(1)} = \bar{r}^{(0)} - \alpha_0 A\bar{r}^{(0)}$. Это дает следующие формулы:

$$\alpha_1 = \frac{(\bar{r}^{(1)}, \bar{r}^{(1)})}{(\bar{r}^{(1)}, A\bar{r}^{(1)})}$$

и

$$\bar{x}^{(2)} = \bar{x}^{(1)} + \alpha_1 \bar{r}^{(1)}.$$

Далее процесс продолжается по формулам:

$$\bar{r}^{(k)} = \bar{f} - A\bar{x}^{(k)} = \bar{r}^{(k-1)} - \alpha_{k-1} A\bar{r}^{(k-1)}, \quad (2.5.58)$$

$$\alpha_k = \frac{(\bar{r}^{(k)}, \bar{r}^{(k)})}{(\bar{r}^{(k)}, A\bar{r}^{(k)})}, \quad (2.5.59)$$

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} + \alpha_k \bar{r}^{(k)} \quad (k=2, 3, \dots). \quad (2.5.60)$$

Заметим, что обычно векторы $\bar{r}^{(k)}$, особенно при большом порядке матрицы системы, удобно вычислять по формуле $\bar{r}^{(k)} = \bar{r}^{(k-1)} - \alpha_{k-1} A\bar{r}^{(k-1)}$. А чтобы из-за ошибок округления так вычисленные векторы $\bar{r}^{(k)}$ через несколько шагов не начали сильно отклоняться от истинных невязок $\bar{f} - A\bar{x}^{(k)}$, их надо время от времени вычислять по формуле $\bar{r}^{(k)} = \bar{f} - A\bar{x}^{(k)}$. В отличие от метода сопряженных градиентов здесь ортогонализация векторов невязок системы $\bar{r}^{(k)}$ не проводится.

Исследуем свойства последовательности векторов $\bar{x}^{(0)}, \bar{x}^{(1)}, \bar{x}^{(2)}, \dots$. Для этой цели нам потребуются две леммы, которые мы ниже и докажем.

Лемма 1. Если a_i — некоторые положительные числа, а γ_i — числа, удовлетворяющие неравенствам $0 < t \leq \gamma_i \leq M$, то справедливо неравенство

$$\frac{\sum_{i=1}^n \gamma_i a_i \sum_{i=1}^n \frac{1}{\gamma_i} a_i}{\left(\sum_{i=1}^n a_i \right)^2} \leq \frac{1}{4} \left[\sqrt{\frac{M}{m}} + \sqrt{\frac{m}{M}} \right]^2. \quad (2.5.61)$$

Доказательство. Введем обозначения

$$\alpha_i = \frac{a_i}{\sum_{i=1}^n a_i} \quad \text{и} \quad \gamma_i = \sqrt{mM} \delta_i.$$

Тогда неравенство (2.5.61) примет вид

$$\sum_{i=1}^n \delta_i \alpha_i \sum_{i=1}^n \frac{1}{\delta_i} \alpha_i \leq \frac{1}{4} \left[\sqrt{\frac{M}{m}} + \sqrt{\frac{m}{M}} \right]^2. \quad (2.5.62)$$

Отметим, что

$$\sum_{i=1}^n \alpha_i = 1 \quad \text{и} \quad \sqrt{\frac{m}{M}} \leq \delta_i \leq \sqrt{\frac{M}{m}}.$$

Так как среднее геометрическое меньше среднего арифметического или равно ему, то будет

$$\sum_{i=1}^n \delta_i \alpha_i \sum_{i=1}^n \frac{1}{\delta_i} \alpha_i \leq \frac{1}{4} \left\{ \sum_{i=1}^n \alpha_i \left(\delta_i + \frac{1}{\delta_i} \right) \right\}^2. \quad (2.5.63)$$

Функция

$$\varphi(\delta) = \delta + \frac{1}{\delta}$$

принимает наибольшее значение на отрезке

$$\left[\sqrt{\frac{m}{M}}, \sqrt{\frac{M}{m}} \right]$$

при

$$\delta = \sqrt{\frac{m}{M}} \quad \text{и} \quad \delta = \sqrt{\frac{M}{m}}.$$

Это значение в обоих случаях равно

$$\sqrt{\frac{m}{M}} + \sqrt{\frac{M}{m}}.$$

Значит,

$$\delta_i + \frac{1}{\delta_i} \leq \sqrt{\frac{m}{M}} + \sqrt{\frac{M}{m}} \quad (2.5.64)$$

при всех $i = 1, 2, \dots, n$.

Теперь из (2.5.63) в силу (2.5.64) получим

$$\begin{aligned} \sum_{i=1}^n \delta_i \alpha_i \sum_{i=1}^n \frac{1}{\delta_i} \alpha_i &\leq \frac{1}{4} \left\{ \sum_{i=1}^n \alpha_i \left[\sqrt{\frac{m}{M}} + \sqrt{\frac{M}{m}} \right] \right\}^2 = \\ &= \frac{1}{4} \left[\sqrt{\frac{m}{M}} + \sqrt{\frac{M}{m}} \right]^2 \left(\sum_{i=1}^n \alpha_i \right)^2 = \frac{1}{4} \left[\sqrt{\frac{m}{M}} + \sqrt{\frac{M}{m}} \right]^2, \end{aligned}$$

ибо $\sum_{i=1}^n \alpha_i = 1$. Лемма доказана.

Введем в рассмотрение понятие функции ошибки, определив ее формулой

$$G(\bar{x}) = (A\bar{e}, \bar{e}), \quad (2.5.65)$$

где $\bar{e} = \bar{x}^{(*)} - \bar{x}$ — вектор ошибки, $\bar{x}^{(*)}$ — точное решение системы $A\bar{x} = \bar{f}$. Имеет место следующая лемма

Лемма 2. Последовательность значений функции ошибки $G(\bar{x}^{(0)})$, $G(\bar{x}^{(1)})$, ..., $G(\bar{x}^{(k)})$, ..., где $\bar{x}^{(k)}$ определяются формулами (2.5.60), стремится к нулю при $k \rightarrow \infty$.
Доказательство. В силу формул (2.5.58) — (2.5.60), (2.5.65) имеем

$$G(\bar{x}^{(k+1)}) - G(\bar{x}^{(k)}) = - \frac{(\bar{r}^{(k)}, \bar{r}^{(k)})^2}{(\bar{r}^{(k)}, A\bar{r}^{(k)})} \quad \text{и} \quad G(\bar{x}^{(k)}) = (A^{-1}\bar{r}^{(k)}, \bar{r}^{(k)}).$$

Значит,

$$\frac{G(\bar{x}^{(k+1)})}{G(\bar{x}^{(k)})} = \frac{1}{G(\bar{x}^{(k)})} \left[G(\bar{x}^{(k)}) - \frac{(\bar{r}^{(k)}, \bar{r}^{(k)})^2}{(\bar{r}^{(k)}, A\bar{r}^{(k)})} \right] = 1 - q_k, \quad (2.5.66)$$

где $q_k = \frac{(\bar{r}^{(k)}, \bar{r}^{(k)})^2}{(\bar{r}^{(k)}, A^{-1}\bar{r}^{(k)}) (\bar{r}^{(k)}, A\bar{r}^{(k)})}$. Оценим снизу величину q_k . Пусть

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n -$$

собственные значения матрицы A и $\bar{u}_1, \bar{u}_2, \dots, \bar{u}_n$ — принадлежащие им собственные векторы, ортогональные друг к другу и нормированные так, что $(\bar{u}_i, \bar{u}_i) = 1$ при $i = 1, 2, \dots, n$. Все $\lambda_i > 0$, ибо A — положительно определенная матрица. Пусть $\lambda_1 \geq m$ и $\lambda_n \leq M$. Разложим вектор $\bar{r}^{(k)}$ по собственным векторам матрицы A :

$$\bar{r}^{(k)} = c_1 \bar{u}_1 + c_2 \bar{u}_2 + \dots + c_n \bar{u}_n. \quad (2.5.67)$$

Так как под $\bar{r}^{(k)}$ мы понимаем ненулевой вектор невязок системы, то в разложении (2.5.67) не все c_i равны нулю. Имеем

$$A\bar{r}^{(k)} = c_1 \lambda_1 \bar{u}_1 + c_2 \lambda_2 \bar{u}_2 + \dots + c_n \lambda_n \bar{u}_n$$

и

$$A^{-1}\bar{r}^{(k)} = c_1 \lambda_1^{-1} \bar{u}_1 + c_2 \lambda_2^{-1} \bar{u}_2 + \dots + c_n \lambda_n^{-1} \bar{u}_n.$$

Следовательно,

$$\begin{aligned}(\bar{r}^{(k)}, \bar{r}^{(k)}) &= \sum_{i=1}^n c_i^2, \quad (\bar{r}^{(k)}, A\bar{r}^{(k)}) = \sum_{i=1}^n \lambda_i c_i^2, \\ (\bar{r}^{(k)}, A^{-1}\bar{r}^{(k)}) &= \sum_{i=1}^n \lambda_i^{-1} c_i^2.\end{aligned}$$

Теперь для q_k получим

$$q_k = \frac{\left(\sum_{i=1}^n c_i^2\right)^2}{\sum_{i=1}^n \lambda_i c_i^2 \sum_{i=1}^n \lambda_i^{-1} c_i^2}.$$

В силу неравенства (2.5.61) из формулы для q_k следует:

$$q_k \geq \frac{4}{\left[\sqrt{\frac{m}{M}} + \sqrt{\frac{M}{m}}\right]^2}.$$

Значит,

$$\frac{G(\bar{x}^{(k+1)})}{G(\bar{x}^{(k)})} = 1 - q_k \leq \left[\frac{M-m}{M+m}\right]^2 < 1.$$

Далее получим

$$G(\bar{x}^{(k+1)}) \leq \left[\frac{M-m}{M+m}\right]^2 G(\bar{x}^{(k)})$$

или

$$\begin{aligned}G(\bar{x}^{(k+1)}) &\leq \left[\frac{M-m}{M+m}\right]^2 G(\bar{x}^{(k)}) \leq \left[\frac{M-m}{M+m}\right]^2 \times \\ &\times \left[\frac{M-m}{M+m}\right]^2 G(\bar{x}^{(k-1)}) \leq \left[\frac{M-m}{M+m}\right]^{2(k+1)} G(\bar{x}^{(0)}).\end{aligned}\quad (2.5.68)$$

Коэффициент $\frac{M-m}{M+m} < 1$, поэтому из (2.5.68) следует, что $G(\bar{x}^{(k+1)}) \rightarrow 0$ при $k \rightarrow \infty$.

Лемма доказана.

Теорема 2. Последовательные приближения $\bar{x}^{(0)}, \bar{x}^{(1)}, \bar{x}^{(2)}, \dots$, построенные по методу скорейшего спуска, сходятся к решению системы $A\bar{x} = \bar{f}$ со скоростью геометрической прогрессии.

Доказательство. Из леммы 2 следует, что $G(\bar{x}^{(k)}) = (A(\bar{x}^{(*)} - \bar{x}^{(k)}), \bar{x}^{(*)} - \bar{x}^{(k)}) \rightarrow 0$ при $k \rightarrow \infty$. А это означает, что $\bar{x}^{(k)} \rightarrow \bar{x}^{(*)}$ при $k \rightarrow \infty$, так как матрица A положительно определенная. Определим теперь скорость сходимости. Имеем

$$G(\bar{x}^{(k)}) = (A\bar{e}^{(k)}, \bar{e}^{(k)}) \geq m|\bar{e}^{(k)}|^2, \quad (2.5.69)$$

где $\bar{e}^{(k)} = \bar{x}^{(*)} - \bar{x}^{(k)}$. Из (2.5.68) и (2.5.69) следует оценка

$$|\bar{\varepsilon}^{(k)}| \leq \sqrt{\frac{G(x^{(k)})}{m}} \leq \sqrt{\frac{G(x^{(0)})}{m}} \left[\frac{M-m}{M+m} \right]^k,$$

означающая, что $|\bar{\varepsilon}^{(k)}|$ стремится к нулю со скоростью геометрической прогрессии. Тем самым утверждение теоремы доказано.

Отметим два свойства приближений $\bar{x}^{(k)}$ метода скорейшего спуска.

1. Невязки двух последовательных приближений ортогональны друг другу.

Действительно, $\bar{r}^{(k+1)} = \bar{r}^{(k)} - \alpha_k A \bar{r}^{(k)}$, откуда $(\bar{r}^{(k+1)}, \bar{r}^{(k)}) = (\bar{r}^{(k)}, \bar{r}^{(k)}) - \alpha_k (A \bar{r}^{(k)}, \bar{r}^{(k)}) = 0$ на основании определения α_k .

2. Каждое последующее приближение ближе к точному решению, чем предыдущее, т. е.

$$\|\bar{x}^{(*)} - \bar{x}^{(k+1)}\| < \|\bar{x}^{(*)} - \bar{x}^{(k)}\|. \quad (2.5.70)$$

Иначе говоря, длина вектора ошибки при переходе к новому приближению строго убывает. Имеем

$$\bar{\varepsilon}^{(k+1)} = \bar{\varepsilon}^{(k)} - \alpha_k \bar{r}^{(k)}.$$

Значит,

$$\begin{aligned} (\bar{\varepsilon}^{(k+1)}, \bar{\varepsilon}^{(k+1)}) &= (\bar{\varepsilon}^{(k)}, \bar{\varepsilon}^{(k)}) - 2\alpha_k (\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) + \alpha_k^2 (\bar{r}^{(k)}, \bar{r}^{(k)}) = \\ &= (\bar{\varepsilon}^{(k)}, \bar{\varepsilon}^{(k)}) - \alpha_k (\bar{\varepsilon}^{(k)}, \bar{\varepsilon}^{(k)}) - \frac{\alpha_k^2}{(\bar{r}^{(k)}, \bar{r}^{(k)})} \left[\frac{(\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) (\bar{r}^{(k)}, \bar{r}^{(k)})}{\alpha_k} - (\bar{r}^{(k)}, \bar{r}^{(k)})^2 \right] = \\ &= (\bar{\varepsilon}^{(k)}, \bar{\varepsilon}^{(k)}) - \alpha_k (\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) - \frac{\alpha_k^2}{(\bar{r}^{(k)}, \bar{r}^{(k)})} \left[(\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) (\bar{r}^{(k)}, A \bar{r}^{(k)}) - (\bar{r}^{(k)}, \bar{r}^{(k)})^2 \right]. \end{aligned}$$

Покажем, что

$$(\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) (\bar{r}^{(k)}, A \bar{r}^{(k)}) - (\bar{r}^{(k)}, \bar{r}^{(k)})^2 \geq 0.$$

Положим $A = B^2$, где B — положительно определенная матрица [2, стр. 113]. Тогда, учитывая, что $\bar{\varepsilon}^{(k)} = \bar{x}^{(*)} - \bar{x}^{(k)} = A^{-1} \bar{f} - \bar{x}^{(k)} = A^{-1} (\bar{f} - A \bar{x}^{(k)}) = A^{-1} \bar{r}^{(k)}$, в силу неравенства Коши — Буняковского получим

$$\begin{aligned} (\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) (\bar{r}^{(k)}, A \bar{r}^{(k)}) - (\bar{r}^{(k)}, \bar{r}^{(k)})^2 &= (\bar{r}^{(k)}, A^{-1} \bar{r}^{(k)}) (\bar{r}^{(k)}, A \bar{r}^{(k)}) - \\ - (\bar{r}^{(k)}, \bar{r}^{(k)})^2 &= (B^{-1} \bar{r}^{(k)}, B^{-1} \bar{r}^{(k)}) (B \bar{r}^{(k)}, B \bar{r}^{(k)}) - (B \bar{r}^{(k)}, B \bar{r}^{(k)})^2 \geq 0. \end{aligned}$$

Таким образом,

$$(\bar{\varepsilon}^{(k+1)}, \bar{\varepsilon}^{(k+1)}) \leq (\bar{\varepsilon}^{(k)}, \bar{\varepsilon}^{(k)}) - \alpha_k (\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) < (\bar{\varepsilon}^{(k)}, \bar{\varepsilon}^{(k)}), \quad (2.5.71)$$

ибо

$$\alpha_k > 0$$

и

$$(\bar{\varepsilon}^{(k)}, \bar{r}^{(k)}) = (\bar{x}^{(*)} - \bar{x}^{(k)}, \bar{f} - A \bar{x}^{(k)}) = (\bar{x}^{(*)} - \bar{x}^{(k)}, A (A^{-1} \bar{f} - \bar{x}^{(k)})) = (\bar{\varepsilon}^{(k)}, A \bar{\varepsilon}^{(k)}) > 0.$$

Следовательно, (2.5.70) следует из (2.5.71).

Метод скорейшего спуска может быть применен и к системам с несимметричной матрицей после умножения слева системы (2.5.25) на матрицу A' . При этом в вычисли-

тельной схеме метода матрица $A'A$ фактически может не вычисляться. Действительно, при указанном умножении

$$A'A\bar{x} = A'\bar{f},$$

и мы должны в качестве невязки взять вектор $\bar{\eta}^{(k)} = A'\bar{r}^{(k)}$, где $\bar{r}^{(k)} = \bar{f} - A\bar{x}^{(k)}$, так, что расчетные формулы будут следующими:

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} + \alpha_k \bar{\eta}^{(k)} = \bar{x}^{(k)} + \alpha_k A'\bar{r}^{(k)}$$

при

$$\alpha_k = \frac{(\bar{\eta}^{(k)}, \bar{\eta}^{(k)})}{(A'A\bar{\eta}^{(k)}, \bar{\eta}^{(k)})} = \frac{(\bar{\eta}^{(k)}, \bar{\eta}^{(k)})}{(A\bar{\eta}^{(k)}, A\bar{\eta}^{(k)})} \quad (k=0, 1, 2, \dots).$$

§ 2.6. СПОСОБЫ ОЦЕНКИ ПОГРЕШНОСТИ ПРИБЛИЖЕННОГО РЕШЕНИЯ СИСТЕМЫ

Как мы уже отмечали, в некоторых методах численного решения систем линейных алгебраических уравнений о точности полученного приближенного решения чаще всего судят по векторам невязок системы. Однако для одного класса матриц малость вектора невязок в некоторой метрике означает и малость компонент вектора погрешностей, для другого класса матриц такой связи может и не быть. Чтобы убедиться в этом, рассмотрим систему

$$A\bar{x} = \bar{f} \quad (2.6.1)$$

и обозначим через $\bar{x}^{(*)}$ ее точное решение, а через \bar{y} — некоторое приближенное решение этой системы. Рассмотрим векторы

$$\bar{\epsilon} = \bar{x}^{(*)} - \bar{y}, \quad \bar{r} = \bar{f} - A\bar{y}, \quad (2.6.2)$$

которые мы будем называть соответственно вектором погрешностей и вектором невязок. Пусть матрица A системы (2.6.1) имеет хотя бы одно очень малое по модулю собственное значение λ , а \bar{z} — соответствующий такому λ собственный вектор этой матрицы. Тогда

$$A(\bar{x}^{(*)} + \bar{z}) = A\bar{x}^{(*)} + A\bar{z} = \bar{f} + \lambda\bar{z}$$

и компоненты вектора $\bar{x}^{(*)} + \bar{z}$ могут отличаться весьма сильно от компонент вектора $\bar{x}^{(*)}$, хотя, в силу малости λ , компоненты вектора $\bar{f} + \lambda\bar{z}$ будут мало отличаться от компонент вектора \bar{f} . В связи с этим необходимо ввести такие соотношения между векторами $\bar{\epsilon}$ и \bar{r} , которые позволяли бы по величине вектора \bar{r} более точно судить о величине вектора $\bar{\epsilon}$. При этом в практике вычислений большее значение имеют не нормы векторов $\bar{\epsilon}$ и \bar{r} , а отношения

$$\frac{\|\bar{\varepsilon}\|}{\|\bar{x}^{(*)}\|} \quad \text{и} \quad \frac{\|\bar{r}\|}{\|\bar{f}\|},$$

являющиеся в некотором смысле «относительными погрешностями». Для количественной характеристики таких отношений, а также и векторов $\bar{\varepsilon}$ и \bar{r} , введем понятия обусловленности систем и матриц.

2.6.1. Обусловленность систем уравнений и матриц

Введем в рассмотрение величину

$$\mu = \sup_{\bar{r}} \left(\frac{\|\bar{\varepsilon}\|}{\|\bar{x}^{(*)}\|} : \frac{\|\bar{r}\|}{\|\bar{f}\|} \right). \quad (2.6.3)$$

Если μ мало, то из (2.6.3) следует, что

$$\|\bar{\varepsilon}\| \leq \mu \frac{\|\bar{x}^{(*)}\|}{\|\bar{f}\|} \|\bar{r}\|,$$

и малость нормы вектора невязок означает малость нормы вектора погрешностей. В этом случае говорят, что система (2.6.1) *хорошо обусловлена*. Если μ велико, то малость нормы $\|\bar{r}\|$ еще не означает малости нормы $\|\bar{\varepsilon}\|$. В этом случае говорят, что система (2.6.1) *плохо обусловлена*. Число μ называют *мерой обусловленности* системы (2.6.1). По аналогии можно ввести и понятие обусловленности матрицы. Из (2.6.2) и определения нормы матрицы имеем

$$\sup_{\bar{r}} \frac{\|\bar{\varepsilon}\|}{\|\bar{r}\|} = \sup_{\bar{r}} \frac{\|\bar{x}^{(*)} - \bar{y}\|}{\|\bar{r}\|} = \sup_{\bar{r}} \frac{\|A^{-1}\bar{r}\|}{\|\bar{r}\|} = \|A^{-1}\|. \quad (2.6.4)$$

Учитывая (2.6.3), из формулы (2.6.4) получим

$$\mu = \frac{\|\bar{f}\|}{\|\bar{x}^{(*)}\|} \|A^{-1}\|. \quad (2.6.5)$$

Будем теперь рассматривать систему (2.6.1) при всевозможных значениях \bar{f} . Тогда решением этой системы будет некоторое множество X — векторов $\bar{x}^{(*)}$, отвечающих соответствующим значениям \bar{f} при одной и той же матрице A . Изучим поведение величины μ , определяемой по формуле (2.6.5), при $\bar{x}^{(*)} \in X$, а именно, вычислим

$$\sup_{\bar{x}^{(*)} \in X} \mu.$$

Имеем

$$\nu = \sup_{\bar{x}^{(*)} \in X} \mu = \sup_{\bar{x}^{(*)} \in X} \frac{\|Ax^{(*)}\|}{\|\bar{x}^{(*)}\|} \cdot \|A^{-1}\| = \|A\| \cdot \|A^{-1}\|. \quad (2.6.6)$$

Назовем ν *числом обусловленности* матрицы A . Из (2.6.6) видно, что если матрица A близка к особенной, то число ν будет для такой матрицы велико. В этом случае говорят, что матрица A *плохо обусловлена*. Если число ν мало, то соответствующую матрицу A называют *хорошо обусловленной*. Как правило, система с плохо (хорошо) обусловленной матрицей A будет плохо (хорошо) обусловленной системой. Значения ν зависят от того, каким образом мы определяем норму матрицы A . Так, например, в случае третьей нормы получим

$$\nu_{III} = \|A\|_{III} \cdot \|A^{-1}\|_{III} = \sqrt{\max \xi_i} \sqrt{\max \pi_i}, \quad (2.6.7)$$

где ξ_i — собственные значения матрицы $A'A$, π_i — собственные значения матрицы $(A^{-1})'A^{-1}$. Так как $(A^{-1})'A^{-1} = (AA')^{-1}$ и матрицы $A'A$, AA' подобны, то собственные значения ξ_i и π_i связаны между собой формулой $\pi_i = \frac{1}{\xi_i}$. Значит, если через λ_n и λ_1 обозначим наибольшее и наименьшее собственные значения матрицы $A'A$, то из (2.6.7) получим

$$\nu_{III} = \|A\|_{III} \cdot \|A^{-1}\|_{III} = \sqrt{\frac{\lambda_n}{\lambda_1}}. \quad (2.6.8)$$

Из (2.6.8) следует, что $\nu_{III} \geq 1$. Это неравенство справедливо для ν при любом выборе нормы матрицы. Действительно,

$$\nu = \|A\| \cdot \|A^{-1}\| \geq \max |\sigma_i| \max |\rho_i| = \max |\sigma_i| \max \left| \frac{1}{\sigma_i} \right| = \frac{\max |\sigma_i|}{\min |\sigma_i|} \geq 1,$$

ибо $\|A\| \geq |\sigma_i|$ и $\|A^{-1}\| \geq |\rho_i|$, где σ_i и ρ_i — собственные значения матриц A и A^{-1} соответственно.

2.6.2. Оценка погрешности $\bar{\varepsilon}$

Такая оценка в сильной степени зависит от того, как изменяется решение системы (2.6.1) при малых изменениях ее коэффициентов и свободных членов. А это означает, что оценка $\bar{\varepsilon}$ зависит от меры и числа обусловленности матрицы системы, т. е. от μ и ν .

Рассмотрим наряду с системой (2.6.1) систему такого вида:

$$B\bar{y} = \bar{g}, \quad (2.6.9)$$

где B и \bar{g} — заданные матрицы и вектор. Предположим, что B и \bar{g} связаны с A и \bar{f} равенствами

$$B = A - CA, \quad \bar{g} = \bar{f} + \bar{\delta}, \quad (2.6.10)$$

где

$$\|C\| \leq q < 1, \quad \|\bar{\delta}\| \leq p.$$

Заметим, что во многих вычислительных алгоритмах приближенное решение системы $A\bar{x} = \bar{f}$ удовлетворяет системам вида (2.6.9), для которых матрица C и вектор $\bar{\delta}$ могут быть вычислены реально. Получим теперь оценку погрешности $\bar{\varepsilon} = \bar{x}^{(*)} - \bar{y}$. Из (2.6.9), (2.6.10) имеем

$$(E - C)A\bar{y} = \bar{g}$$

или

$$\begin{aligned} A\bar{y} &= (E - C)^{-1}\bar{g} = (E + C + C^2 + \dots)(\bar{f} + \bar{\delta}) = \\ &= \bar{f} + (C + C^2 + \dots)\bar{f} + (E + C + C^2 + \dots)\bar{\delta}. \end{aligned}$$

Сравнивая эту формулу с формулой $\bar{r} = \bar{f} - A\bar{y}$, замечаем, что вектор $\bar{r} = -[(C + C^2 + \dots)\bar{f} + (E + C + C^2 + \dots)\bar{\delta}]$ можно рассматривать как невязку при приближенном решении \bar{y} системы (2.6.1). Таким образом, в силу определения μ и ν , имеем

$$\begin{aligned} \frac{\|\bar{\varepsilon}\|}{\|\bar{x}^{(*)}\|} &= \frac{\|\bar{x}^{(*)} - \bar{y}\|}{\|\bar{x}^{(*)}\|} \leq \mu \frac{\|\bar{r}\|}{\|\bar{f}\|} = \mu \frac{\|-(C + C^2 + \dots)\bar{f} - (E + C + C^2 + \dots)\bar{\delta}\|}{\|\bar{f}\|} \leq \\ &\leq \mu \left\{ \frac{q}{1-q} + \frac{1}{1-q} \cdot \frac{\|\bar{\delta}\|}{\|\bar{f}\|} \right\} \leq \nu \left\{ \frac{q}{1-q} + \frac{1}{1-q} \cdot \frac{p}{\|\bar{f}\|} \right\}. \quad (2.6.11) \end{aligned}$$

Отсюда видно, что «относительная погрешность» $\frac{\|\bar{\varepsilon}\|}{\|\bar{x}^{(*)}\|}$ тем меньше, чем меньше число обусловленности ν . Малость такой погрешности сильно зависит также от того, сколь сильно уклоняются матрица B и вектор \bar{g} соответственно от матрицы A и вектора \bar{f} , т. е. зависит от малости величин q и p . Из (2.6.11) можно получить далее оценку для $\|\bar{\varepsilon}\|$, которой

можно будет пользоваться в практике вычислений. Обозначим

$$l(p, q) = \frac{q}{1-q} + \frac{1}{1-q} \cdot \frac{p}{\|\bar{f}\|}.$$

Учитывая, что $\|\bar{y} + \bar{x}^{(*)} - \bar{y}\| \leq \|\bar{y}\| + \|\bar{x}^{(*)} - \bar{y}\| \leq \|\bar{y}\| + \|\bar{\varepsilon}\|$, из (2.6.11) получим

$$\|\bar{\varepsilon}\| \leq \|\bar{y}\| \cdot \frac{v \cdot l(p, q)}{1 - v \cdot l(p, q)} \quad (2.6.12)$$

при условии, что $1 - v \cdot l(p, q) > 0$. Реально нам известны матрица B и вектор \bar{g} , а не A и \bar{f} , поэтому вместо оценки (2.6.12) можно рассматривать следующую оценку:

$$\|\bar{\varepsilon}\| \leq \|\bar{y}\| \frac{v^{(*)} \cdot l(p, q^{(*)})}{1 - l(p, q^{(*)})}, \quad (2.6.13)$$

где $v^{(*)}$ — число обусловленности матрицы B , $q^{(*)} = \|DB^{-1}\|$, $D = B - A$

$$\text{и } l(p, q^{(*)}) = \frac{q^{(*)}}{1 - q^{(*)}} + \frac{1}{1 - q^{(*)}} \cdot \frac{p}{\|\bar{g}\|}.$$

Большое число методов решения системы вида (2.6.1) основано на преобразовании матрицы A к некоторому простейшему виду (например, к диагональному, треугольному и т. д.). Чаще всего такое преобразование выполняется путем умножения матрицы A слева на некоторую невырожденную матрицу M . В связи с этим выясним, какой класс матриц M при указанном преобразовании не меняет числа обусловленности матрицы A , т. е. определим матрицы M , для которых имеет место равенство

$$v(MA) = v(A). \quad (2.6.14)$$

Для любой невырожденной матрицы M верно неравенство

$$v(MA) \leq v(M)v(A), \quad (2.6.15)$$

ибо

$$\|MA\| \leq \|M\| \cdot \|A\| \quad \text{и} \quad \|A^{-1} \cdot M^{-1}\| \leq \|A^{-1}\| \cdot \|M^{-1}\|.$$

С другой стороны,

$$v(MA) \geq \frac{1}{v(M)} v(A), \quad (2.6.16)$$

так как на основании неравенства (2.6.15)

$$v(A) \leq v(M^{-1})v(MA) = v(M)v(MA).$$

Аналогично получим

$$\nu(M \cdot A) \geq \frac{1}{\nu(A)} \nu(M). \quad (2.6.17)$$

Таким образом,

$$\max \left(\frac{1}{\nu(M)} \nu(A), \frac{1}{\nu(A)} \nu(M) \right) \leq \nu(M \cdot A) \leq \nu(M) \nu(A). \quad (2.6.18)$$

Из (2.6.18) видно, что если при некоторой заданной матрице M правая граница неравенства достигается для какой-либо матрицы A и $\nu(M)$ велико, то число обусловленности $\nu(M \cdot A)$ может стать очень большим.

Если положить $\nu(M) = 1$, то из (2.6.18) сразу будет следовать (2.6.14). Значит, все невырожденные матрицы, у которых число обусловленности равно единице ($\nu(M) = 1$), не меняют числа обусловленности матрицы A , т. е. $\nu(M \cdot A) = \nu(A)$.

Отметим, что указанным свойством в случае третьей нормы, очевидно, обладают ортогональные и унитарные матрицы.

Литература

1. Березин И. С., Жидков Н. П. Методы вычислений, т. 1. М., 1966.
2. Воеводин В. В. Численные методы алгебры (теория и алгоритмы). М., 1966.
3. Ланцош К. Практические методы прикладного анализа. М., 1961.
4. Фаддеев Д. К., Фаддеева В. Н. Вычислительные методы линейной алгебры. М., 1963.
5. Форсайт Дж., Молер К. Численное решение систем линейных алгебраических уравнений. М., 1969.
6. Уилкинсон Дж. Х. Алгебраическая проблема собственных значений. М., 1970.
7. Хаусхолдер А. С. Основы численного анализа. М., 1956.
8. Forsythe G. Решение линейных алгебраических уравнений может быть интересным. Bull. Amer. Math. Soc., 59, № 64, 1953.

Глава 3

ВЫЧИСЛЕНИЕ СОБСТВЕННЫХ ЗНАЧЕНИЙ И СОБСТВЕННЫХ ВЕКТОРОВ МАТРИЦ

§ 3.1. О СОДЕРЖАНИИ ЗАДАЧИ

В предыдущей главе мы ознакомились лишь с одной из основных групп вычислительных задач линейной алгебры — с задачами численного нахождения решения системы линейных алгебраических уравнений. Сейчас мы рассмотрим другую важную группу таких задач, порождаемую так называемой проблемой собственных значений.

Как мы уже отмечали, *собственным значением* (или *характеристическим числом*) квадратной матрицы A называется такое число λ , что для некоторого ненулевого вектора \bar{x} имеет место равенство

$$A\bar{x} = \lambda\bar{x}. \quad (3.1.1)$$

Любой ненулевой вектор \bar{x} , удовлетворяющий этому равенству, называется *собственным вектором* матрицы A , соответствующим (или принадлежащим) собственному значению λ . Очевидно, что все собственные векторы матрицы определены с точностью до числового множителя.

Уже в предыдущей главе мы имели возможность убедиться, насколько ценной бывает информация о собственных значениях матрицы. Например, скорость и сам факт сходимости процесса простых итераций, применяемого для приближенного решения системы линейных алгебраических уравнений вида

$$\bar{x} = B\bar{x} + \bar{b},$$

существенным образом зависят от величины максимального по модулю собственного значения матрицы B . Задача нахождения собственных значений и собственных векторов матрицы важна не только как вспомогательная. Многие научно-технические задачи (особенно задачи физики, механики, астрономии) приводят к проблеме отыскания нетривиального решения однородной системы линейных алгебраических уравнений вида (3.1.1) и тех значений числового параметра λ , при которых такое решение существует. Во всех явлениях неустойчивых колебаний и вибраций проблема собственных значений играет очень важную роль, так как частота колебаний определяется собственными значениями некоторой матрицы,

а форму этих колебаний указывают собственные векторы этой матрицы. Анализ собственных значений матриц является важной темой научно-технических исследований.

Условием существования у однородной системы (3.1.1) ненулевого решения (для наглядности запишем эту систему в виде $(A - \lambda E)\bar{x} = \bar{0}$) является требование

$$|A - \lambda E| = \begin{vmatrix} a_{11} - \lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \lambda \end{vmatrix} = 0.$$

Это уравнение обычно называют *вековым* (или *характеристическим*) *уравнением* матрицы A . Такие уравнения часто встречаются в приложениях. Левая часть векового уравнения

$$|A - \lambda E| = (-1)^n (\lambda^n - p_1 \lambda^{n-1} - p_2 \lambda^{n-2} - \dots - p_n)$$

носит название *характеристического полинома* матрицы A . Старший коэффициент этого полинома равен $(-1)^n$. Иногда вместо характеристического полинома рассматривают полином, отличающийся от характеристического множителем $(-1)^n$. Этот полином

$$P(\lambda) = \lambda^n - p_1 \lambda^{n-1} - p_2 \lambda^{n-2} - \dots - p_n$$

обычно называют *собственным многочленом* матрицы. Собственные значения матрицы являются корнями собственного многочлена. Совокупность всех собственных значений $\lambda_1, \lambda_2, \dots, \lambda_n$ матрицы A , где каждое собственное значение выписано столько раз, какова его кратность как корня собственного многочлена, называется *спектром* этой матрицы. Собственными же векторами матрицы A являются нетривиальные решения однородной системы (3.1.1), в которой вместо λ подставлены собственные значения λ_i матрицы. В том случае, когда для данного собственного значения система (3.1.1) имеет несколько линейно независимых решений, этому собственному значению принадлежит несколько собственных векторов. Отметим, что в случае вещественной матрицы комплексному собственному значению соответствуют собственные векторы, координаты которых также будут комплексными числами. Вектор, координаты которого комплексно сопряжены с координатами собственного вектора вещественной матрицы, также будет собственным вектором данной матрицы, соответствующим комплексно сопряженному собственному значению ее. В этом легко убедиться, если в равенстве $Ax = \lambda x$ заменить все числа комплексно сопряженными.

Задачу вычисления собственных значений и собственных векторов матрицы A можно разбить на три естественных этапа:

- 1) построение собственного многочлена $P(\lambda)$ матрицы;
- 2) решение уравнения $P(\lambda) = 0$ и нахождение собственных значений λ_i ($i = 1, 2, \dots, n$) матрицы;
- 3) отыскание нетривиальных решений однородных систем

$$(A - \lambda_i E) \bar{x} = \bar{0} \quad (i = 1, 2, \dots, n),$$

т. е. нахождение собственных векторов матрицы.

Как мы увидим в дальнейшем, иногда можно вычислять собственные значения и принадлежащие им собственные векторы матрицы, минуя этап построения собственного многочлена этой матрицы. Этого удается достигнуть при помощи различных косвенных соображений, использующих те или иные свойства собственных значений и собственных векторов матрицы.

Каждый из трех отмеченных этапов решения проблемы собственных значений представляет собой достаточно сложную вычислительную задачу.

В самом деле, построение собственного многочлена $P(\lambda)$, например, связано с развертыванием определителя

$$\begin{vmatrix} a_{11} - \lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \lambda \end{vmatrix} = (-1)^n (\lambda^n - p_1 \lambda^{n-1} - p_2 \lambda^{n-2} - \dots - p_n) = (-1)^n P(\lambda), \quad (3.1.2)$$

что представляет собой значительные технические трудности. Основное затруднение вызвано тем обстоятельством, что λ входит в каждую строку и в каждый столбец определителя. В общем же случае, как известно из алгебры, коэффициенты p_i собственного многочлена $P(\lambda)$ представляют собой взятые со знаком $(-1)^{i-1}$ суммы всех главных миноров (т. е. миноров, симметрично расположенных относительно главной диагонали) порядка i определителя матрицы A . Число таких миноров для каждого i равно числу сочетаний из n по i . Значит, непосредственное вычисление коэффициентов собственного многочлена $P(\lambda)$ квадратной матрицы порядка n связано с вычислением

$$C_n^1 + C_n^2 + \dots + C_n^n = 2^n - 1$$

определителей различных порядков. Для матриц достаточно высоких порядков последняя задача сопряжена с большими затратами вычислительного труда.

Полезно отметить, что в силу известной теоремы Виета, дающей связь корней многочлена с его коэффициентами, можно записать следующие равенства:

$$\lambda_1 + \lambda_2 + \dots + \lambda_n = p_1,$$

$$\lambda_1 \cdot \lambda_2 \dots \lambda_n = (-1)^{n-1} p_n.$$

Так как в силу равенства (3.1.2)

$$p_1 = a_{11} + a_{22} + \dots + a_{nn},$$

$$p_n = (-1)^{n-1} |A|,$$

то

$$\lambda_1 + \lambda_2 + \dots + \lambda_n = a_{11} + a_{22} + \dots + a_{nn} = \text{Sp } A,$$

$$\lambda_1 \cdot \lambda_2 \dots \lambda_n = |A|.$$

Таким образом, сумма всех собственных значений матрицы совпадает с ее следом, а произведение их равно значению определителя этой матрицы. В частности, отсюда следует, что матрица A тогда и только тогда имеет хотя бы одно собственное значение, равное нулю, когда $|A| = 0$, т. е. если она особенная.

Трудности в непосредственном осуществлении второго и третьего этапов решения проблемы собственных значений, т. е. трудности, связанные с решением алгебраических уравнений высоких степеней, и трудности в нахождении нетривиальных решений систем однородных линейных алгебраических уравнений, также значительны. После ознакомления с содержанием первых двух глав книги уже нетрудно оценить объем вычислительной работы, необходимый для непосредственного осуществления этих этапов рассматриваемой проблемы.

К настоящему времени создано немало специальных вычислительных приемов, упрощающих численное нахождение собственных значений и собственных векторов матрицы. Все эти методы, как и в случае проблемы численного решения системы линейных алгебраических уравнений, можно разделить на точные и итерационные методы. К первой группе относятся методы, по которым сначала строят собственный многочлен матрицы (т. е. вычисляют его коэффициенты p_1, p_2, \dots, p_n), затем, находя его корни, получают собственные значения матрицы и уже по ним находят соответствующие собственные векторы. При этом во многих случаях, используя промежуточные результаты вычислений, получают собственные векторы матрицы, принадлежащие вычисленным собственным значениям, не прибегая к решению указанных выше систем однородных линейных алгебраических уравнений. Методы этой группы получили название точных методов в связи с тем обстоятельством, что в случае точного задания (рациональными числами) элементов матрицы и при точном

(по правилам действий над обыкновенными дробями) проведении вычислений такие методы приводят к точным значениям коэффициентов собственного многочлена, а координаты собственных векторов при этом оказываются выраженными через соответствующие собственные значения.

В методах второй группы собственные значения матрицы определяются непосредственно, без обращения к собственному многочлену, при этом обычно одновременно вычисляются и соответствующие собственные векторы. Вычислительные схемы таких методов носят итерационный характер. В них используется многократное умножение матрицы на вектор. Схемы этого типа обычно приводят к последовательности векторов, имеющей своим пределом собственный вектор, и к числовой последовательности, предел которой является соответствующим собственным значением. При этом ход итерационного процесса существенным образом зависит от характера канонической формы Жордана для данной матрицы, а также от наличия у матрицы вещественных или комплексных собственных значений. Сам факт сходимости этого процесса и ее скорость определяются величиной отношения модулей различных соседних собственных значений.

Как правило, итерационные методы позволяют с достаточной точностью определить лишь первые (наибольшие по модулю, например) собственные значения и соответствующие им собственные векторы. Поэтому методы этой группы чаще всего применяются к решению так называемой *частичной проблемы собственных значений*; т. е. их чаще используют лишь для отыскания одного или нескольких собственных значений матрицы и соответствующих собственных векторов. Точные же методы позволяют решать также и *полную проблему собственных значений*, т. е. дают возможность находить все собственные значения матрицы и все принадлежащие им собственные векторы. Полная проблема собственных значений в некоторых случаях может быть решена также и специальными итерационными методами. Эти методы, конечно, более трудоемки, чем точные методы и чем итерационные методы решения частичной проблемы собственных значений. Их практическое использование стало возможным лишь с появлением быстродействующих вычислительных машин. Однако перед точными методами решения полной проблемы собственных значений итерационные методы имеют одно несомненное преимущество, связанное с возможностью нахождения всех собственных значений без предварительного построения собственного многочлена матрицы. Это особенно важно в связи с тем, что ошибки в вычислении коэффициентов собственного многочлена могут сильно сказываться на точности определения его корней, т. е. на точности нахождения собственных значений исходной матрицы (и соответствующих им собственных векторов). Кроме того, большим достоинством итерационных методов перед точными является простота и единообразие производимых действий, что особенно ценно при использовании быстродействующих вычислительных машин.

Полная и частичная проблемы собственных значений сильно различаются как по методам их решения, так и по области приложений. Так как решение полной проблемы собственных значений даже в случае матриц не очень высокого порядка обычно связано с очень большим объемом вычислительного труда, то возможность решения частичной проблемы собственных значений другими методами, минуя вычислительные трудности решения полной проблемы, является очень ценной для практики.

Изложение вычислительных методов решения проблемы собственных значений мы начнем с рассмотрения группы точных методов, при этом, если противное не оговорено особо, мы будем иметь в виду лишь матрицы с вещественными элементами.

§ 3.2. МЕТОД А. Н. КРЫЛОВА

В начале тридцатых годов нашего столетия А. Н. Крыловым был предложен достаточно удобный метод нахождения собственных значений и собственных векторов матриц. Сообщение об этом методе положило начало большому циклу работ, посвященных приведению векового урав-

$$|A - \lambda E| = \begin{vmatrix} a_{11} - \lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \lambda \end{vmatrix} = 0 \quad (3.2.1)$$

к полиномиальному виду

$$(-1)^n(\lambda^n - p_1\lambda^{n-1} - p_2\lambda^{n-2} - \dots - p_n) = 0. \quad (3.2.2)$$

К настоящему времени вычислительная схема метода значительно улучшена, однако основная идея метода не претерпела больших изменений. Для иллюстрации ее А. Н. Крылов вводит в рассмотрение каноническую систему однородных обыкновенных дифференциальных уравнений первого порядка с постоянными коэффициентами

$$\begin{aligned} y_1' &= a_{11}y_1 + a_{12}y_2 + \dots + a_{1n}y_n, \\ y_2' &= a_{21}y_1 + a_{22}y_2 + \dots + a_{2n}y_n, \\ &\vdots \\ y_n' &= a_{n1}y_1 + a_{n2}y_2 + \dots + a_{nn}y_n, \end{aligned}$$

связанную с исходной матрицей A . Характеристическое уравнение этой системы имеет вид (3.2.1). Корни характеристического уравнения системы являются собственными значениями матрицы A . Если эту систему

уравнений первого порядка удастся свести к одному дифференциальному уравнению порядка n с постоянными коэффициентами

$$y^{(n)} = p_1 y^{(n-1)} + p_2 y^{(n-2)} + \dots + p_n y,$$

то по виду этого уравнения легко записать его характеристическое уравнение

$$\lambda^n - p_1 \lambda^{n-1} - p_2 \lambda^{n-2} - \dots - p_n = 0,$$

корни которого должны совпадать с корнями уравнения (3.2.1). Итак, выполнив преобразование введенной системы обыкновенных дифференциальных уравнений первого порядка к одному уравнению порядка n , мы построим математический образ, по виду которого непосредственно записывается вековое уравнение исходной матрицы в полиномиальном виде (3.2.2). Во многих случаях такой прием оказался не только возможным, но и достаточно удобным для вычислений.

А. Н. Крылов указал и на возможность алгебраической интерпретации этой идеи, хотя сам разработкой ее не занимался. Мы сейчас остановимся именно на построении такого алгебраического образа, по виду которого можно будет непосредственно записать собственный многочлен $P(\lambda)$ матрицы A или его делитель, при этом, оказывается, результаты промежуточных алгебраических преобразований могут быть использованы и для вычисления собственных векторов матрицы.

Прежде чем приступить к рассмотрению этой алгебраической интерпретации метода Крылова, мы приведем здесь некоторые сведения из высшей алгебры, необходимые нам для изложения.

3.2.1. Некоторые сведения из алгебры

Назовем многочлен

$$f(\lambda) = a_0 \lambda^m + a_1 \lambda^{m-1} + \dots + a_m$$

аннулирующим многочленом для квадратной матрицы A , если

$$f(A) \equiv a_0 A^m + a_1 A^{m-1} + \dots + a_{m-1} A + a_m E = 0.$$

Нулевой многочлен является аннулирующим для любой матрицы. Будем рассматривать только приведенные (со старшим коэффициентом, равным единице) аннулирующие многочлены. Для каждой матрицы множество таких многочленов не пусто. Действительно, в алгебре матриц хорошо известна теорема Гамильтона — Кели, утверждающая, что *если $P(\lambda)$ есть собственный многочлен матрицы A , то $P(A) = 0$* , т. е., условно говоря, *матрица является корнем своего собственного многочлена*. Таким образом, любая квадратная матрица порядка n имеет аннулирующий многочлен n -й степени. Очевидно, что этот многочлен не единственный,

так как если многочлен $P(\lambda)$ является аннулирующим для матрицы A , то этим свойством обладает и всякий многочлен, делящийся на $P(\lambda)$. Среди всего множества многочленов $f(\lambda)$, аннулирующих для данной матрицы A , особо выделяют многочлен $\psi(\lambda)$ наименьшей степени. Такой многочлен называется *минимальным многочленом матрицы*. Укажем на некоторые почти очевидные свойства этого многочлена.

1. Если $f(A) = 0$, то многочлен $f(\lambda)$ делится нацело на минимальный многочлен $\psi(\lambda)$ матрицы A .

В самом деле, пусть

$$f(\lambda) = \psi(\lambda) Q(\lambda) + r(\lambda),$$

где многочлен $r(\lambda)$ имеет степень меньшую, чем многочлен $\psi(\lambda)$. Покажем, что это возможно лишь в случае $r(\lambda) \equiv 0$. Действительно, подстановка в последнее равенство A вместо λ приводит нас к результату

$$r(A) = 0,$$

который возможен лишь в случае $r(\lambda) \equiv 0$, ибо многочлен $\psi(\lambda)$ имеет наименьшую степень среди всех многочленов, аннулирующих для матрицы A .

2. Все корни минимального многочлена матрицы являются собственными значениями этой матрицы.

Это свойство является прямым следствием свойства 1.

Можно показать также, что корнями минимального многочлена матрицы служат все различные между собой собственные значения матрицы.

3. Минимальный многочлен матрицы единствен.

Действительно, если $\psi_1(\lambda)$ и $\psi_2(\lambda)$ есть два минимальных многочлена матрицы A , то многочлен меньшей степени

$$Q(\lambda) = \psi_1(\lambda) - \psi_2(\lambda)$$

будет аннулирующим многочленом для этой матрицы, а это может быть лишь в случае $Q(\lambda) \equiv 0$, т. е. при $\psi_1(\lambda) \equiv \psi_2(\lambda)$, так как приведенные многочлены $\psi_1(\lambda)$ и $\psi_2(\lambda)$ имеют наименьшую степень среди всех многочленов, для которых матрица A является корнем.

Пусть, далее, наряду с квадратной матрицей A имеется некоторый вектор $\bar{c} \neq 0$, согласованный по размерности с матрицей A . Рассмотрим множество приведенных многочленов $g(\lambda)$ таких, что

$$g(A)\bar{c} = \bar{0}.$$

Очевидно, что множество таких многочленов $g(\lambda)$ включает в себя рассмотренное выше множество многочленов $f(\lambda)$, для которых матрица A является корнем. Этому множеству принадлежат, в частности, собственный и минимальный многочлены матрицы, но ему могут принадлежать

и такие многочлены, для которых условие $g(A)=0$ не выполняется. Среди многочленов этого множества также особо выделяют многочлен $\varphi(\lambda)$ наименьшей степени, который обычно называют *минимальным аннулирующим вектор \bar{c} многочленом матрицы A* . Как и в случае минимального многочлена $\psi(\lambda)$ матрицы A , можно проверить, что введенный нами минимальный аннулирующий вектор \bar{c} многочлен $\varphi(\lambda)$ матрицы A обладает свойствами, аналогичными свойствам минимального многочлена матрицы.

1. Если $g(A)\bar{c}=\bar{0}$, то многочлен $g(\lambda)$ нацело делится на минимальный аннулирующий вектор \bar{c} многочлен $\varphi(\lambda)$ матрицы A .

2. Все корни многочлена $\varphi(\lambda)$ являются собственными значениями матрицы A . Корни многочлена $\varphi(\lambda)$ дают, вообще говоря, только часть различных собственных значений матрицы.

3. Минимальный аннулирующий вектор \bar{c} многочлен матрицы A единствен.

Приведенные сведения из алгебры будут необходимы при рассмотрении алгебраической интерпретации метода А. Н. Крылова, к описанию которой мы и переходим.

3.2.2. Нахождение собственных значений матрицы

Рассмотрим произвольный вектор $\bar{c}^{(0)} \neq \bar{0}$, согласованный по размерности с исходной квадратной матрицей A . Очень часто в качестве вектора $\bar{c}^{(0)}$ берут, например, вектор $(1, 0, 0, \dots, 0)'$. По этому вектору $\bar{c}^{(0)}$ будем составлять последовательность векторов $\bar{c}^{(1)}=A\bar{c}^{(0)}$, $\bar{c}^{(2)}=A\bar{c}^{(1)}=A^2\bar{c}^{(0)}$, $\bar{c}^{(3)}=A^3\bar{c}^{(0)}$ и т. д. до тех пор, пока не встретим первый вектор (например, вектор $\bar{c}^{(m)}=A^m\bar{c}^{(0)}$), который будет являться линейной комбинацией предыдущих линейно независимых векторов, т. е. пока не будет справедливо следующее равенство:

$$q_1\bar{c}^{(m-1)}+q_2\bar{c}^{(m-2)}+\dots+q_m\bar{c}^{(0)}=\bar{c}^{(m)} \quad \left(\sum_{i=1}^m q_i^2 > 0 \right).$$

Очевидно, что $m \leq n$, где n — размерность вектора $\bar{c}^{(0)}$. Для того чтобы практически определить число m и найти коэффициенты q_1, q_2, \dots, q_m соответствующей линейной комбинации, можно поступить следующим образом.

Запишем предельно возможную ($m=n$) линейную комбинацию

$$q_1\bar{c}^{(n-1)}+q_2\bar{c}^{(n-2)}+\dots+q_n\bar{c}^{(0)}=\bar{c}^{(n)}$$

покоординатно

$$\begin{aligned} q_1 c_1^{(m-1)} + q_2 c_1^{(m-2)} + \dots + q_m c_1^{(0)} &= c_1^{(m)}, \\ q_1 c_2^{(m-1)} + q_2 c_2^{(m-2)} + \dots + q_m c_2^{(0)} &= c_2^{(m)}, \\ &\vdots \\ q_1 c_n^{(m-1)} + q_2 c_n^{(m-2)} + \dots + q_m c_n^{(0)} &= c_n^{(m)}. \end{aligned}$$

и выбрав (например, по методу Гаусса) из этих n линейных алгебраических уравнений m линейно независимых, мы найдем коэффициенты q_1, q_2, \dots, q_m разыскиваемой линейной комбинации.

Оказывается, построенная линейная комбинация и будет тем алгебраическим образом, по виду которого можно непосредственно записать либо собственный многочлен матрицы (при $m=n$), либо его делитель (при $m < n$).

Рассмотрим сначала случай, когда $m=n$. В этом случае, оказывается, коэффициенты q_1, q_2, \dots, q_n линейной комбинации

$$q_1 \bar{c}^{(n-1)} + q_2 \bar{c}^{(n-2)} + \dots + q_n \bar{c}^{(0)} = \bar{c}^{(n)}$$

равны соответствующим коэффициентам p_1, p_2, \dots, p_n собственного многочлена

$$P(\lambda) = \lambda^n - p_1\lambda^{n-1} - p_2\lambda^{n-2} - \dots - p_n,$$

т. е. $q_i = p_i \quad (i = 1, 2, \dots, n)$.

Действительно, на основании теоремы Гамильтона — Кели

$$P(A) \equiv A^n - p_1 A^{n-1} - p_2 A^{n-2} - \dots - p_n E = 0.$$

Умножая это равенство на вектор $\bar{s}^{(0)}$ и принимая во внимание, что

$$A^i \bar{c}^{(0)} = \bar{c}^{(i)} \quad (i=1, 2, \dots, n),$$

получим

$$p_1 \bar{c}^{(n-1)} + p_2 \bar{c}^{(n-2)} + \dots + p_n \bar{c}^{(0)} = \bar{c}^{(n)}.$$

С другой стороны,

$$q_1 \bar{c}^{(n-1)} + q_2 \bar{c}^{(n-2)} + \dots + q_n \bar{c}^{(0)} = \bar{c}^{(n)}.$$

Значит,

$$(p_1 - q_1) \bar{c}^{(n-1)} + (p_2 - q_2) \bar{c}^{(n-2)} + \dots + (p_n - q_n) \bar{c}^{(0)} = \bar{0}.$$

Так как векторы $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(n-1)}$ линейно независимы, то последнее равенство возможно лишь в случае, если

$$p_i = q_i \quad (i = 1, 2, \dots, n).$$

Таким образом, в случае $m=n$ по виду построенной нами линейной комбинации можно непосредственно записать собственный многочлен $P(\lambda)$ матрицы A . Решая уравнение $P(\lambda)=0$ (например, методом Ньютона), мы найдем все собственные значения этой матрицы.

В случае $m < n$ построенная линейная комбинация имеет вид

$$q_1 \bar{c}^{(m-1)} + q_2 \bar{c}^{(m-2)} + \dots + q_m \bar{c}^{(0)} = \bar{c}^{(m)}.$$

Если учесть, что $\bar{c}^{(i)} = A^i \bar{c}^{(0)}$ ($i=1, 2, \dots, m$), то последнее равенство можно переписать в виде

$$(A^m - q_1 A^{m-1} - q_2 A^{m-2} - \dots - q_m E) \bar{c}^{(0)} = \bar{0}$$

или

$$\varphi(A) \bar{c}^{(0)} = \bar{0},$$

где

$$\varphi(\lambda) = \lambda^m - q_1 \lambda^{m-1} - q_2 \lambda^{m-2} - \dots - q_m.$$

Следовательно, найдя коэффициенты q_1, q_2, \dots, q_m искомой линейной комбинации, мы тем самым построим многочлен $\varphi(\lambda)$, который будет являться минимальным аннулирующим вектор $\bar{c}^{(0)}$ многочленом матрицы A (если бы существовал многочлен $g(\lambda)$, удовлетворяющий условию $g(A) \bar{c}^{(0)} = \bar{0}$ и имеющий степень меньшую, чем степень многочлена $\varphi(\lambda)$, то это противоречило бы условию линейной независимости векторов $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(m-1)}$).

Таким образом, в случае $m < n$ по виду построенной линейной комбинации мы сможем записать не сам собственный многочлен $P(\lambda)$ матрицы A , а лишь его делитель $\varphi(\lambda)$. Решив уравнение

$$\varphi(\lambda) = 0,$$

мы найдем лишь часть собственных значений этой матрицы. Изменяя исходный вектор $\bar{c}^{(0)}$, можно на этом пути найти и недостающие собственные значения.

3.2.3. Вычисление собственных векторов матрицы

После того как собственное значение λ_i матрицы A вычислено, задача нахождения принадлежащих ему собственных векторов этой матрицы сводится, вообще говоря, к решению следующей однородной системы линейных алгебраических уравнений:

$$(A - \lambda_i E) \bar{x} = \bar{0}.$$

Но часто промежуточные результаты вычислений при нахождении собственных значений матрицы могут быть с успехом использованы и для вычисления соответствующих собственных векторов. Это, как правило, позволяет сократить затраты вычислительного труда при решении последнего этапа проблемы собственных значений. Такая возможность, в частности, представляется и описанным выше алгоритмом метода Крылова.

Пусть известен корень λ_i минимального аннулирующего вектор $\bar{c}^{(0)}$ многочлена

$$\varphi(\lambda) = \lambda^m - q_1\lambda^{m-1} - q_2\lambda^{m-2} - \dots - q_m$$

матрицы A (все последующие рассуждения имеют место как для регулярного случая $m=n$, так и для особого случая $m < n$). Собственный вектор $\bar{x}^{(i)}$ матрицы A , принадлежащий этому собственному значению, будем искать в виде линейной комбинации линейно независимых векторов $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(m-1)}$, построенных при нахождении многочлена $\varphi(\lambda)$:

$$\bar{x}^{(i)} = \beta_{i1}\bar{c}^{(m-1)} + \beta_{i2}\bar{c}^{(m-2)} + \dots + \beta_{im}\bar{c}^{(0)}. \quad (3.2.3)$$

Коэффициенты β_{ij} ($j=1, 2, \dots, m$) надлежит выбрать так, чтобы удовлетворить условию

$$A\bar{x}^{(i)} = \lambda_i\bar{x}^{(i)}.$$

Умножим линейную комбинацию (3.2.3) на матрицу A , учитывая равенства $\bar{c}^{(j)} = A\bar{c}^{(j-1)}$ ($j=1, 2, \dots, m$) и требование $A\bar{x}^{(i)} = \lambda_i\bar{x}^{(i)}$:

$$\lambda_i(\beta_{i1}\bar{c}^{(m-1)} + \beta_{i2}\bar{c}^{(m-2)} + \dots + \beta_{im}\bar{c}^{(0)}) = \beta_{i1}\bar{c}^{(m)} + \beta_{i2}\bar{c}^{(m-1)} + \dots + \beta_{im}\bar{c}^{(1)}. \quad (3.2.4)$$

Если же, кроме того, учесть, что $\varphi(A)\bar{c}^{(0)} = \bar{0}$, т. е. что

$$\bar{c}^{(m)} = q_1\bar{c}^{(m-1)} + q_2\bar{c}^{(m-2)} + \dots + q_m\bar{c}^{(0)},$$

равенство (3.2.4) можно переписать в виде

$$\begin{aligned} & \lambda_i(\beta_{i1}\bar{c}^{(m-1)} + \beta_{i2}\bar{c}^{(m-2)} + \dots + \beta_{im}\bar{c}^{(0)}) = \\ & = \beta_{i1}(q_1\bar{c}^{(m-1)} + q_2\bar{c}^{(m-2)} + \dots + q_m\bar{c}^{(0)}) + \beta_{i2}\bar{c}^{(m-1)} + \beta_{i3}\bar{c}^{(m-2)} + \dots + \beta_{im}\bar{c}^{(1)} \end{aligned}$$

или

$$\begin{aligned} & (q_m\beta_{i1} - \lambda_i\beta_{im})\bar{c}^{(0)} + (q_{m-1}\beta_{i1} + \beta_{im} - \lambda_i\beta_{im-1})\bar{c}^{(1)} + \\ & + (q_{m-2}\beta_{i1} + \beta_{im-1} - \lambda_i\beta_{im-2})\bar{c}^{(2)} + \dots + (q_1\beta_{i1} + \beta_{i2} - \lambda_i\beta_{i1})\bar{c}^{(m-1)} = \bar{0}. \end{aligned}$$

§ 3.3. МЕТОД А. М. ДАНИЛЕВСКОГО

Достаточно простой и экономичный способ решения проблемы собственных значений был предложен в конце тридцатых годов этого столетия А. М. Данилевским. Этот метод основан на известном из линейной алгебры факте о том, что преобразование подобия $S^{-1}AS$ не изменяет характеристического полинома матрицы A . Действительно,

$$|S^{-1}AS - \lambda E| = |S^{-1}AS - \lambda S^{-1}ES| = |S^{-1}| \cdot |A - \lambda E| \cdot |S| = |A - \lambda E|.$$

Поэтому, удачно подобрав преобразование подобия, можно надеяться получить матрицу, собственный многочлен которой выписывается непосредственно по виду ее. А. М. Данилевский предложил приводить исходную матрицу A преобразованием подобия $S^{-1}AS$ к так называемой канонической форме Фробениуса

$$\Phi = \begin{bmatrix} p_1 & p_2 & p_3 & \dots & p_{n-1} & p_n \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & 0 \end{bmatrix},$$

характеристический полином которой легко записать. В самом деле, разлагая определитель $|\Phi - \lambda E|$ последовательно по элементам первого столбца, будем иметь

$$\begin{aligned} |\Phi - \lambda E| &= \begin{vmatrix} p_1 - \lambda & p_2 & p_3 & \dots & p_{n-1} & p_n \\ 1 & -\lambda & 0 & \dots & 0 & 0 \\ 0 & 1 & -\lambda & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & -\lambda \end{vmatrix} = (p_1 - \lambda)(-\lambda)^{n-1} - \\ &- \begin{vmatrix} p_2 & p_3 & \dots & p_{n-1} & p_n \\ 1 & -\lambda & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & & 1 & -\lambda \end{vmatrix} = (p_1 - \lambda)(-\lambda)^{n-1} - p_2(-\lambda)^{n-2} + \\ &+ \begin{vmatrix} p_3 & p_4 & \dots & p_{n-1} & p_n \\ 1 & -\lambda & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & & 1 & -\lambda \end{vmatrix} = \\ &= (p_1 - \lambda)(-\lambda)^{n-1} - p_2(-\lambda)^{n-2} + p_3(-\lambda)^{n-3} - \dots + (-1)^{n+1}p_n = \\ &= (-1)^n(\lambda^n - p_1\lambda^{n-1} - p_2\lambda^{n-2} - \dots - p_n) = (-1)^nP(\lambda). \end{aligned}$$

Таким образом, элементы p_1, p_2, \dots, p_n первой строки матрицы Фробениуса являются соответствующими коэффициентами ее собственного многочлена, а значит, и собственного многочлена исходной матрицы A , связанной с матрицей Φ преобразованием подобия $\Phi = S^{-1}AS$.

Решая уравнение $P(\lambda) = 0$, мы найдем интересующие нас собственные значения матрицы A . Кроме того, оказывается, неособенная матрица S , с помощью которой было построено нужное нам преобразование подобия, может быть использована при нахождении собственных векторов матрицы A .

Основная задача, таким образом, сводится к разысканию нужной нам матрицы S . А. М. Данилевский предложил строить эту матрицу и тем самым осуществлять переход от матрицы A к матрице Φ последовательно с помощью $n-1$ преобразований подобия, переводящих строки матрицы A , начиная с последней, в соответствующие строки матрицы Φ . Рассмотрим эти преобразования подробнее.

3.3.1. Построение собственного многочлена матрицы

В зависимости от элементов матрицы A в методе Данилевского можно встретиться с двумя возможными случаями: регулярным и нерегулярным. Рассмотрим сначала регулярный случай.

Предположим, что элемент a_{nn-1} матрицы A отличен от нуля. Тогда, разделив $(n-1)$ -й столбец матрицы A на этот элемент и вычитая этот столбец из i -го столбца матрицы, домножив его предварительно на элемент a_{ni} (для всех $i=1, 2, \dots, n-2, n$), мы приведем последнюю строку матрицы к форме Фробениуса. Непосредственно проверяется, что такое преобразование равносильно умножению матрицы A справа на матрицу

$$M_{n-1} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ -\frac{a_{n1}}{a_{nn-1}} & -\frac{a_{n2}}{a_{nn-1}} & \dots & -\frac{a_{nn-2}}{a_{nn-1}} & \frac{1}{a_{nn-1}} & -\frac{a_{nn}}{a_{nn-1}} \\ 0 & 0 & \dots & 0 & 0 & 1 \end{bmatrix}.$$

В результате такого умножения последняя строка матрицы принимает нужный вид, однако преобразование AM_{n-1} не будет, вообще говоря, преобразованием подобия для матрицы A . Исправить этот недостаток можно умножением полученной матрицы слева на матрицу M_{n-1}^{-1} , которая существует, так как $|M_{n-1}| = \frac{1}{a_{nn-1}} \neq 0$. Непосредственно убеждаемся, что матрица M_{n-1}^{-1} имеет следующий вид:

$$M_{n-1}^{-1} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn-1} & a_{nn} \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}.$$

Очевидно, что преобразование $M_{n-1}^{-1}AM_{n-1}$ не изменяет последней строки матрицы AM_{n-1} .

Таким образом, после выполнения первого шага метода Данилевского мы получим матрицу следующего вида:

$$M_{n-1}^{-1}AM_{n-1} = A^{(1)} = \begin{bmatrix} a_{11}^{(1)} & \dots & a_{12}^{(1)} & \dots & a_{1n-1}^{(1)} & a_{1n}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & \dots & a_{2n-1}^{(1)} & a_{2n}^{(1)} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ a_{n-11}^{(1)} & a_{n-12}^{(1)} & \dots & a_{n-1n-1}^{(1)} & a_{n-1n}^{(1)} \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix}.$$

Заметим, что матрицы M_{n-1}^{-1} и M_{n-1} , умножением на которые соответственно слева и справа мы переходим от матрицы A к матрице $A^{(1)}$, выписываются непосредственно по виду матрицы A .

Предположим далее, что и элемент $a_{n-2n-1}^{(1)}$ матрицы $A^{(1)}$ отличен от нуля (имеем в виду регулярный случай). Тогда второй шаг метода Данилевского совершенно аналогичен первому и состоит в приведении второй снизу строки матрицы $A^{(1)}$ к форме Фробениуса (при сохранении неизменной первой снизу строки). Результат таких преобразований можно записать в виде

$$\begin{aligned} M_{n-2}^{-1}M_{n-1}^{-1}AM_{n-1}M_{n-2} &= M_{n-2}^{-1}A^{(1)}M_{n-2} = \\ &= A^{(2)} = \begin{bmatrix} a_{11}^{(2)} & \dots & a_{1n-2}^{(2)} & a_{1n-1}^{(2)} & a_{1n}^{(2)} \\ \dots & \dots & \dots & \dots & \dots \\ a_{n-21}^{(2)} & \dots & a_{n-2n-2}^{(2)} & a_{n-2n-1}^{(2)} & a_{n-2n}^{(2)} \\ 0 & \dots & 1 & 0 & 0 \\ 0 & \dots & 0 & 1 & 0 \end{bmatrix}, \end{aligned}$$

где

$$M_{n-2} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \frac{a_{n-11}^{(1)}}{a_{n-1n-2}^{(1)}} & \frac{a_{n-12}^{(1)}}{a_{n-1n-2}^{(1)}} & \dots & \frac{a_{n-1n-3}^{(1)}}{a_{n-1n-2}^{(1)}} & \frac{1}{a_{n-1n-2}^{(1)}} & \frac{a_{n-1n-1}^{(1)}}{a_{n-1n-2}^{(1)}} & \frac{a_{n-1n}^{(1)}}{a_{n-1n-2}^{(1)}} \\ 0 & 0 & \dots & 0 & 0 & 1 & 0 \\ 0 & 0 & \dots & 0 & 0 & 0 & 1 \end{bmatrix},$$

$$M_{n-2}^{-1} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ a_{n-11}^{(1)} & a_{n-12}^{(1)} & \dots & a_{n-1n-1}^{(1)} & a_{n-1n}^{(1)} \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}.$$

Закон построения матриц M_{n-2}^{-1} и M_{n-2} по виду матрицы $A^{(1)}$, как видим, вполне аналогичен соответствующему правилу построения на предыдущем шаге метода матриц M_{n-1}^{-1} и M_{n-1} по виду матрицы $A = A^{(0)}$. Эта же закономерность сохраняется и на последующих шагах метода.

Итак, если $a_{nn-1}^{(1)} \neq 0$, $a_{n-1n-2}^{(1)} \neq 0$, $a_{n-2n-3}^{(2)} \neq 0$, ..., $a_{21}^{(n-2)} \neq 0$, то после $n-1$ шагов метода Данилевского будем иметь

$$M_1^{-1} M_2^{-1} \dots M_{n-1}^{-1} A M_{n-1} M_{n-2} \dots M_1 = A^{(n-1)} =$$

$$= \begin{bmatrix} a_{11}^{(n-1)} & a_{12}^{(n-1)} & \dots & a_{1n-1}^{(n-1)} & a_{1n}^{(n-1)} \\ 1 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} = \begin{bmatrix} p_1 & p_2 & \dots & p_{n-1} & p_n \\ 1 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} =$$

$$= \Phi = S^{-1} A S.$$

Тем самым исходная матрица A посредством преобразования подобия с неособенной матрицей $S = M_{n-1} M_{n-2} \dots M_1$ будет приведена к канонической форме Фробениуса, непосредственно по виду первой строки которой записывается собственный многочлен

$$P(\lambda) = \lambda^n - p_1 \lambda^{n-1} - p_2 \lambda^{n-2} - \dots - p_n.$$

Рассмотрим далее нерегулярный случай. Будем считать, что процесс последовательного приведения строк исходной матрицы A к виду Фробениуса по методу Данилевского доведен до строки номера k , т. е. выполнено $n-k$ шагов метода, но при этом оказалось, что элемент $a_{kk-k}^{(n-k)}$ матрицы $A^{(n-k)}$ равен нулю. Следующий $(n-k+1)$ -й шаг метода изложенным выше способом осуществлен быть не может. В зависимости от того, есть ли среди элементов k -й строки матрицы $A^{(n-k)}$, стоящих левее элемента $a_{kk-k}^{(n-k)}=0$, отличные от нуля или таковых нет, дальнейшее продолжение процесса возможно, например, по двум следующим вариантам.

Предположим сначала, что имеет место первая из двух оговоренных выше возможностей, т. е. в строке номера k левее элемента $a_{k, k-1}^{(n-k)} = 0$ есть элемент, отличный от нуля. Пусть, к примеру, этот элемент стоит в i -м ($i < k-1$) столбце матрицы $A^{(n-k)}$. Тогда дальнейшее продолжение процесса может быть сведено к регулярному случаю. Для этого, оказывается, достаточно в матрице $A^{(n-k)}$ поменять местами столбцы с номерами i и $k-1$, а также строки с такими номерами. Непосредственно легко проверить, что такое преобразование может быть записано в виде

$$TA^{(n-k)}T,$$

где

$$T = \begin{bmatrix} 1 & 0 & & & & & & & & & \\ 0 & 1 & & & & & & & & & \\ & & \ddots & & & & & & & & \\ & & & 1 & & & & & & & \\ & & & & \ddots & & & & & & \\ & & & & & 0 & \dots & \dots & \dots & \dots & \dots \\ & & & & & & 1 & \dots & \dots & \dots & \dots \\ & & & & & & & \ddots & & & \\ & & & & & & & & 1 & \dots & \dots \\ & & & & & & & & & 0 & \dots & \dots \\ & & & & & & & & & & 1 & \dots & \dots \\ & & & & & & & & & & & \ddots & & \\ & & & & & & & & & & & & 1 & 0 \\ & & & & & & & & & & & & 0 & 1 \end{bmatrix} \quad \begin{matrix} (i) \\ \\ \\ \\ \\ (i) \\ \\ \\ \\ (k-1) \end{matrix}$$

Легко проверяется также, что преобразование $TA^{(n-k)}T$ есть преобразование подобия для матрицы $A^{(n-k)}$. В самом деле, поскольку после двойной перестановки строк или столбцов мы получаем исходную матрицу,

то $T^2 = E$, т. е. $T = T^{-1}$. Значит, преобразование $TA^{(n-k)}T$ есть подобное преобразование матрицы $A^{(n-k)}$.

Проведя такое преобразование (дополнительные затраты труда на это невелики), мы сможем следующий шаг метода Данилевского выполнять, как и в регулярном случае.

Рассмотрим сейчас вторую возможность, которая может представиться в нерегулярном случае, т. е. предположим, что

$$a_{k1}^{(n-k)} = a_{k2}^{(n-k)} = \dots = a_{kk-1}^{(n-k)} = 0.$$

Матрица $A^{(n-k)}$ в этом случае имеет вид

$$A^{(n-k)} = \begin{bmatrix} a_{11}^{(n-k)} & \dots & a_{1k-1}^{(n-k)} & a_{1k}^{(n-k)} & \dots & a_{1n-1}^{(n-k)} & a_{1n}^{(n-k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ a_{k-11}^{(n-k)} & \dots & a_{k-1k-1}^{(n-k)} & a_{k-1k}^{(n-k)} & \dots & a_{k-1n-1}^{(n-k)} & a_{k-1n}^{(n-k)} \\ 0 & \dots & 0 & a_{kk}^{(n-k)} & \dots & a_{kn-1}^{(n-k)} & a_{kn}^{(n-k)} \\ 0 & \dots & 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & 0 & \dots & 1 & 0 \end{bmatrix} = \begin{bmatrix} B^{(n-k)} & C^{(n-k)} \\ 0 & \Phi^{(n-k)} \end{bmatrix},$$

где

$$B^{(n-k)} = \begin{bmatrix} a_{11}^{(n-k)} & a_{12}^{(n-k)} & \dots & a_{1k-1}^{(n-k)} \\ a_{21}^{(n-k)} & a_{22}^{(n-k)} & \dots & a_{2k-1}^{(n-k)} \\ \dots & \dots & \dots & \dots \\ a_{k-11}^{(n-k)} & a_{k-12}^{(n-k)} & \dots & a_{k-1k-1}^{(n-k)} \end{bmatrix},$$

$$\Phi^{(n-k)} = \begin{bmatrix} a_{kk}^{(n-k)} & \dots & a_{kn-1}^{(n-k)} & a_{kn}^{(n-k)} \\ 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & 1 & 0 \end{bmatrix}.$$

Тогда

$$|A^{(n-k)} - \lambda E| = |B^{(n-k)} - \lambda E_{k-1}| \cdot |\Phi^{(n-k)} - \lambda E_{n-k+1}|. (*)$$

*) Это равенство является простым следствием теоремы Лапласа (см., например, А. Г. Курош. Курс высшей алгебры. М., 1962, гл. I, § 6). Индексами снизу в правой части равенства обозначены порядки единичных матриц.

Так как матрица $\Phi^{(n-k)}$ есть матрица Фробениуса, то ее характеристический многочлен выписывается непосредственно по виду первой строки. Значит, для нахождения многочлена $|A^{(n-k)} - \lambda E|$ достаточно привести к канонической форме Фробениуса лишь квадратную матрицу $B^{(n-k)}$ порядка $k-1 < n$. Таким образом, в этом случае задача построения собственного многочлена матрицы даже упрощается.

Нетрудно простым подсчетом необходимых арифметических операций убедиться в том, что метод Данилевского является одним из самых экономичных среди известных методов построения собственного многочлена матрицы. Однако, как и почти все точные методы, он очень чувствителен к ошибкам в результатах промежуточных вычислений. Известным уже нам простым приемом можно несколько повысить надежность вычислений в методе Данилевского, если на $(n-k+1)$ -м шаге метода на место элемента $a_{hh-1}^{(n-k)}$ ставить с помощью преобразования подобия наибольший по модулю среди элементов матрицы $A^{(n-k)}$, стоящих выше или левее элемента $a_{hh-1}^{(n-k)}$. Для контроля вычислений при этом полезно сравнивать полученное значение коэффициента p_1 со следом матрицы.

3.3.2. Вычисление собственных векторов матрицы

Если найдены собственные значения λ_i ($i=1, 2, \dots, n$) матрицы A и известна неособенная матрица S , преобразование подобия с помощью которой приводит исходную матрицу к канонической форме Фробениуса, то в методе Данилевского, как и в случае метода Крылова, при нахождении собственных векторов матрицы A можно обойтись и без решения систем однородных линейных алгебраических уравнений

$$A\bar{x} = \lambda_i \bar{x} \quad (i=1, 2, \dots, n).$$

Результаты промежуточных вычислений при нахождении собственных значений матрицы здесь также могут быть использованы и для вычисления собственных векторов этой матрицы.

Как мы уже отмечали, матрицы, связанные преобразованием подобия, имеют одинаковые спектры. Собственные же векторы этих матриц, принадлежащие одним и тем же собственным значениям, будут, вообще говоря, различны. Но между ними существует связь, а именно: *если вектор \bar{x} есть собственный вектор матрицы A , принадлежащий собственному значению λ , а вектор \bar{y} — собственный вектор подобной ей матрицы $\Phi = S^{-1}AS$, принадлежащий тому же собственному значению λ , то вектор $S\bar{y}$ также будет собственным вектором матрицы A , соответствующим собственному значению λ .*

Действительно, так как $\Phi\bar{y} = \lambda\bar{y}$ и $\Phi = S^{-1}AS$, то

$$S^{-1}AS\bar{y} = \lambda\bar{y}.$$

Умножая это равенство слева на матрицу S , получаем утверждаемое:

$$AS\bar{y} = \lambda S\bar{y}.$$

Таким образом, собственные векторы исходной матрицы A легко находятся по соответствующим собственным векторам ее канонической формы Фробениуса. Проблема же нахождения собственных векторов матрицы Фробениуса решается просто. Действительно, если λ — известное значение матрицы Φ , то

$$\Phi\bar{y} = \lambda\bar{y}$$

или

$$\begin{bmatrix} p_1 & p_2 & \dots & p_{n-1} & p_n \\ 1 & 0 & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \lambda \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

Запишем это векторное равенство покоординатно:

$$\begin{aligned} p_1 y_1 + p_2 y_2 + \dots + p_n y_n &= \lambda y_1, \\ y_1 &= \lambda y_2, \\ \cdot & \\ y_2 &= \lambda y_3, \\ \cdot & \cdot \cdot \cdot \\ y_{n-1} &= \lambda y_n. \end{aligned}$$

Принимая во внимание, что собственный вектор матрицы определен с точностью до постоянного множителя, положим $y_n = 1$. Тогда из предыдущих равенств можно последовательно найти остальные координаты вектора \bar{y} :

$$y_{n-1} = \lambda, y_{n-2} = \lambda^2, \dots, y_1 = \lambda^{n-1}.$$

Равенство же

$$p_1 y_1 + p_2 y_2 + \dots + p_n y_n = \lambda y_1$$

при этом принимает тривиальный вид

$$P(\lambda) \equiv \lambda^n - p_1 \lambda^{n-1} - p_2 \lambda^{n-2} - \dots - p_n = 0.$$

Его можно использовать для контроля вычислений.

Таким образом, вектор $(\lambda^{n-1}, \lambda^{n-2}, \dots, \lambda, 1)'$ является собственным вектором матрицы Фробениуса, соответствующим собственному значению λ .

Итак, зная матрицу S , нетрудно решить и задачу нахождения собственных векторов исходной матрицы A . Если нахождение собственных

значений производилось по методу Данилевского, то матрица S непосредственно выписывается в регулярном случае метода и в первом варианте нерегулярного случая (когда перестановками соответствующих строк и столбцов вычисления сводятся к регулярному случаю). Например, в регулярном случае

$$S = M_{n-1}M_{n-2} \dots M_1.$$

Так как матрицы M_i ($i = 1, 2, \dots, n-1$) только одной строкой отличаются от единичной, то вектор

$$\bar{x} = S\bar{y} = M_{n-1}M_{n-2} \dots M_1\bar{y}$$

удобнее строить, не находя предварительно произведения $S = M_{n-1}M_{n-2} \dots M_1$, а производя умножения вектора \bar{y} последовательно на матрицы M_1, M_2, \dots, M_{n-1} . При этом от умножения на матрицу M_i будет, очевидно, изменяться лишь i -я координата вектора.

При втором варианте в нерегулярном случае метода Данилевского, когда последовательное приведение (снизу) строк данной матрицы к форме Фробениуса провести до конца не удастся, использовать описанный выше прием, значительно облегчающий задачу вычисления собственных векторов исходной матрицы, естественно, нельзя.

§ 3.4. ДРУГИЕ МЕТОДЫ ПОЛУЧЕНИЯ СОБСТВЕННОГО МНОГОЧЛЕНА МАТРИЦЫ

Выше мы рассмотрели два точных метода решения полной проблемы собственных значений: метод Крылова и метод Данилевского. К настоящему времени известно также большое число других методов, позволяющих находить собственный многочлен матрицы. Не ставя перед собой задачи дать полное и подробное изложение этих методов, мы ограничимся здесь лишь кратким обзором некоторых из них.

3.4.1. Интерполяционный метод

Как мы уже отмечали ранее, задача построения характеристического многочлена трудна тем, что требует непосредственного разворачивания определителя. Метод интерполяции позволяет заменить трудоемкую задачу разворачивания определителя $|A - \lambda E|$ более простой задачей вычисления значений этого определителя при фиксированных значениях переменной λ . Идея метода основана на хорошо известном из алгебры факте о том, что алгебраический многочлен степени n вполне определяется своими значениями в $n+1$ точках. Так как старший коэффициент интересующего нас характеристического многочлена равен $(-1)^n$, то для восстановления остальных его n коэффициентов достаточно подсчитать n значений определителя $|A - \lambda E|$.

Зададим n любых различных значений λ_i ($i=1, 2, \dots, n$) переменной λ и подсчитаем значения $D_i = |A - \lambda_i E|$ ($i=1, 2, \dots, n$) рассматриваемого определителя. Эту задачу можно решить одним из численных методов, рассмотренных в предыдущей главе. Тогда для определения коэффициентов p_1, p_2, \dots, p_n собственного многочлена $P(\lambda)$ мы получим следующую систему n линейных алгебраических уравнений:

$$\begin{aligned} (-1)^n (\lambda_1^n - p_1 \lambda_1^{n-1} - p_2 \lambda_1^{n-2} - \dots - p_n) &= D_1, \\ (-1)^n (\lambda_2^n - p_1 \lambda_2^{n-1} - p_2 \lambda_2^{n-2} - \dots - p_n) &= D_2, \\ &\vdots \\ (-1)^n (\lambda_n^n - p_1 \lambda_n^{n-1} - p_2 \lambda_n^{n-2} - \dots - p_n) &= D_n. \end{aligned}$$

Определитель этой системы с точностью до знака совпадает с определителем Вандермонда

$$\begin{vmatrix} \lambda_1^{n-1} & \lambda_1^{n-2} & \dots & \lambda_1 & 1 \\ \lambda_2^{n-1} & \lambda_2^{n-2} & \dots & \lambda_2 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \lambda_n^{n-1} & \lambda_n^{n-2} & \dots & \lambda_n & 1 \end{vmatrix}.$$

Так как значения λ_i ($i=1, 2, \dots, n$) попарно различны, то этот определитель отличен от нуля. Следовательно, коэффициенты собственного многочлена $P(\lambda)$ могут быть найдены и при этом единственным образом.

Описанный метод называют интерполяционным, так как задача точного или приближенного восстановления функции по нескольким известным ее значениям является простейшей задачей теории интерполирования, которая будет подробно изложена в следующей главе этой книги. Там же будет указан и ряд способов построения интерполяционного многочлена, которые позволят после подсчета значений $D_i = |A - \lambda_i E|$ ($i = 1, 2, \dots, n$) находить коэффициенты собственного многочлена $P(\lambda)$ матрицы A , минуя задачу решения выпянной выше системы линейных алгебраических уравнений.

Интерполяционный метод построения собственного многочлена матрицы, хотя и значительно упрощает задачу непосредственного развертывания определителя $|A - \lambda E|$, все же остается достаточно громоздким, так как требует вычисления n значений определителя. Этот метод удобен, если матрица A имеет невысокий порядок и если легко можно подобрать такие значения переменной λ , для которых определитель $|A - \lambda E|$ вычисляется просто. Кроме того, метод интерполяции не позволяет как-нибудь упростить задачу нахождения собственных векторов матрицы, в то время как методы, например, Крылова и Данилевского значительно облегчают решение этой задачи. Однако интерполяционный метод важен и интересен прежде всего широкой областью его применимости, а также тем, что

он позволяет решать и более общие задачи. В самом деле, так как для этого метода специальный вид определителя, дающего характеристический многочлен, не имеет значения, то он, очевидно, с успехом может быть применен к задаче разворачивания определителя

$$\begin{vmatrix} P_{11}(\lambda) & P_{12}(\lambda) & \dots & P_{1n}(\lambda) \\ P_{21}(\lambda) & P_{22}(\lambda) & \dots & P_{2n}(\lambda) \\ \dots & \dots & \dots & \dots \\ P_{n1}(\lambda) & P_{n2}(\lambda) & \dots & P_{nn}(\lambda) \end{vmatrix},$$

где $P_{ij}(\lambda)$ ($i, j=1, 2, \dots, n$) — известные алгебраические многочлены переменной λ .

3.4.2. Метод Лаверье

Этот метод является хронологически одним из первых методов, предложенных для решения рассматриваемой проблемы. Несмотря на большой объем работы, обусловленный вычислительной схемой метода, он давно получил признание как один из универсальных и простых по логике алгоритма методов построения собственного многочлена матрицы.

Идея метода Лаверье основана на использовании хорошо известных из алгебры формул Ньютона

$$kp_k = S_k - p_1 S_{k-1} - p_2 S_{k-2} - \dots - p_{k-1} S_1 \quad (k=1, 2, \dots, n), \quad (3.4.1)$$

связывающих коэффициенты p_1, p_2, \dots, p_n собственного многочлена

$$P(\lambda) = \lambda^n - p_1 \lambda^{n-1} - p_2 \lambda^{n-2} - \dots - p_n$$

матрицы A с симметрическими функциями

$$S_k = \sum_{i=1}^n \lambda_i^k \quad (k=1, 2, \dots, n)$$

его корней, т. е. собственных значений $\lambda_1, \lambda_2, \dots, \lambda_n$ этой матрицы.

Если значения S_k известны, то формулы Ньютона (3.4.1) позволяют последовательно вычислять коэффициенты собственного многочлена матрицы:

$$\begin{aligned} p_1 &= S_1, \\ p_2 &= \frac{1}{2} (S_2 - S_1 p_1), \\ p_3 &= \frac{1}{3} (S_3 - S_2 p_1 - S_1 p_2), \\ &\dots \\ p_n &= \frac{1}{n} (S_n - S_{n-1} p_1 - S_{n-2} p_2 - \dots - S_1 p_{n-1}). \end{aligned}$$

Величины же S_k принципиально нетрудно подсчитать по исходной матрице A . В самом деле, так как собственные значения матрицы A^k есть $\lambda_1^k, \lambda_2^k, \dots, \lambda_n^k$, то

$$S_k = \lambda_1^k + \lambda_2^k + \dots + \lambda_n^k = \text{Sp } A^k = \sum_{i=1}^n a_{ii}^{(k)},$$

где через $a_{ij}^{(k)}$ ($i, j = 1, 2, \dots, n$) обозначены элементы матрицы A^k .

Таким образом, для нахождения собственного многочлена $P(\lambda)$ матрицы A нужно по этой матрице составить матрицы A^2, A^3, \dots, A^n , подсчитать следы $S_k = \text{Sp } A^k$ ($k = 1, 2, \dots, n$) этих матриц и по формулам Ньютона получить коэффициенты p_1, p_2, \dots, p_n искомого многочлена.

Если учесть, что для вычисления следа матрицы нужно знать не все ее элементы, а лишь диагональные, то можно ограничиться составлением матриц A^k лишь при $k = 2, 3, \dots, m$, где $m = \left\lfloor \frac{n+1}{2} \right\rfloor$. Следы же матриц $A^{m+1}, A^{m+2}, \dots, A^n$ теперь уже можно находить, минуя вычисление недиагональных элементов этих матриц. Это позволяет значительно сократить объем вычислений.

Но и при учете сделанного замечания метод Леверье остается очень трудоемким, так как он связан с многократным умножением матриц.

3.4.3. Метод Д. К. Фаддеева

Интересное видоизменение метода Леверье было предложено Д. К. Фаддеевым. Оно не только позволяет вычислять коэффициенты собственного многочлена матрицы, но и дает возможность эффективно находить матрицу, обратную данной, а также может быть использовано и для получения собственных векторов исходной матрицы.

Предлагается вместо последовательности матриц A, A^2, \dots, A^n находить другую матричную последовательность A_1, A_2, \dots, A_n , построенную следующим образом:

$$\begin{aligned} A_1 &= A, & \text{Sp } A_1 &= q_1, & B_1 &= A_1 - q_1 E, \\ A_2 &= AB_1, & \frac{\text{Sp } A_2}{2} &= q_2, & B_2 &= A_2 - q_2 E, \\ &\dots & \dots & \dots & \dots & \dots \\ A_{n-1} &= AB_{n-2}, & \frac{\text{Sp } A_{n-1}}{n-1} &= q_{n-1}, & B_{n-1} &= A_{n-1} - q_{n-1} E, \\ A_n &= AB_{n-1}, & \frac{\text{Sp } A_n}{n} &= q_n, & B_n &= A_n - q_n E. \end{aligned}$$

При этом, оказывается, будут справедливы следующие утверждения:

- 1) $q_i = p_i$ ($i = 1, 2, \dots, n$),
- 2) матрица B_n есть нулевая матрица,

3) если матрица A — неособенная, то

$$A^{-1} = \frac{B_{n-1}}{p_n}.$$

Проверим сначала первое утверждение. Для доказательства равенств $q_i = p_i$ ($i=1, 2, \dots, n$) применим метод математической индукции. При $i=1$ справедливость утверждаемого очевидна:

$$p_1 = \text{Sp } A = \text{Sp } A_1 = q_1.$$

Предположив, что выполняются равенства $q_i = p_i$ для всех $i=1, 2, \dots, k-1$, докажем, что $q_k = p_k$. Так как по построению

$$A_k = A^k - q_1 A^{k-1} - q_2 A^{k-2} - \dots - q_{k-1} A,$$

а по предположению $q_i = p_i$ для $i=1, 2, \dots, k-1$, то

$$A_k = A^k - p_1 A^{k-1} - p_2 A^{k-2} - \dots - p_{k-1} A.$$

Следовательно,

$$\begin{aligned} kq_k &= \text{Sp } A_k = \text{Sp } A^k - p_1 \text{Sp } A^{k-1} - p_2 \text{Sp } A^{k-2} - \dots - p_{k-1} \text{Sp } A = \\ &= S_k - p_1 S_{k-1} - p_2 S_{k-2} - \dots - p_{k-1} S_1. \end{aligned}$$

Но в силу формул Ньютона (3.4.1)

$$S_k - p_1 S_{k-1} - p_2 S_{k-2} - \dots - p_{k-1} S_1 = kp_k.$$

Значит, $kq_k = kp_k$, что и доказывает справедливость первого утверждения.

Второе утверждение также легко доказать, если воспользоваться теоремой Гамильтона — Кели:

$$B_n = A_n - q_n E = A^n - p_1 A^{n-1} - p_2 A^{n-2} - \dots - p_{n-1} A - p_n E = 0.$$

Проверим, наконец, последнее утверждение. Так как по только что доказанному $A_n = p_n E$, а по построению $A_n = AB_{n-1}$, то

$$AB_{n-1} = p_n E$$

или

$$A^{-1} = \frac{1}{p_n} B_{n-1},$$

что и требовалось доказать.

Можно показать, что в случае особенной матрицы A матрица $C = (-1)^{n-1} B_{n-1}$ будет союзной с матрицей A , т. е.

$$C = (A_{ij})' \quad (i, j=1, 2, \dots, n),$$

где через A_{ij} обозначено алгебраическое дополнение элемента a_{ij} в определителе матрицы A .

Заметим, что доказанное ранее равенство $A_n = p_n E$ может быть использовано и для контроля вычислений: об их точности можно судить по отклонению матрицы A_n от скалярной.

Метод Фаддеева позволяет также эффективно находить и собственные векторы матрицы: для этих целей используются промежуточные результаты вычислений, производимых при построении собственного многочлена матрицы.

Рассмотрим матрицу

$$Q(\lambda) = \lambda^{n-1}E + \lambda^{n-2}B_1 + \lambda^{n-3}B_2 + \dots + \lambda B_{n-2} + B_{n-1}.$$

Можно показать, что если все собственные значения $\lambda_1, \lambda_2, \dots, \lambda_n$ исходной матрицы A различны, то матрицы $Q(\lambda_i)$ ($i=1, 2, \dots, n$) — ненулевые. В этом случае, оказывается, любой ненулевой столбец матрицы $Q(\lambda_i)$ может быть принят в качестве собственного вектора матрицы A , соответствующего собственному значению λ_i .

В самом деле,

$$\begin{aligned} (\lambda_i E - A) Q(\lambda_i) &= (\lambda_i E - A) (\lambda_i^{n-1} E + \lambda_i^{n-2} B_1 + \lambda_i^{n-3} B_2 + \dots + \lambda_i B_{n-2} + B_{n-1}) = \\ &= \lambda_i^n E + \lambda_i^{n-1} (B_1 - A) + \lambda_i^{n-2} (B_2 - AB_1) + \dots + \lambda_i (B_{n-1} - AB_{n-2}) - AB_{n-1} = \\ &= (\lambda_i^n - p_1 \lambda_i^{n-1} - p_2 \lambda_i^{n-2} - \dots - p_n) E = 0, \end{aligned}$$

так как по построению $B_k - AB_{k-1} = -p_k E$ ($k=1, 2, \dots, n$), а λ_i есть корень собственного многочлена.

Из полученного равенства

$$(\lambda_i E - A) Q(\lambda_i) = 0$$

следует, что

$$(\lambda_i E - A) \bar{x} = \bar{0}$$

или

$$A \bar{x} = \lambda_i \bar{x},$$

где \bar{x} — любой столбец матрицы $Q(\lambda_i)$.

Таким образом, любой ненулевой столбец матрицы $Q(\lambda_i)$ может быть принят в качестве собственного вектора матрицы A , принадлежащего собственному значению λ_i .

При нахождении собственных векторов матрицы A таким способом нет необходимости, конечно, строить всю матрицу $Q(\lambda_i)$, а достаточно для каждого λ_i ($i=1, 2, \dots, n$) ограничиться вычислением лишь одного ее столбца.

В случае кратных собственных значений задача нахождения соответствующих собственных векторов усложняется. Наряду с матрицей $Q(\lambda)$ здесь может понадобиться привлекать к рассмотрению также матрицы, полученные дифференцированием ее по λ .

3.4.4. Метод окаймления

Идея окаймления, с которой мы встречались уже в проблеме нахождения решения системы линейных алгебраических уравнений (см. п. 2.3.3) может быть полезной и в проблеме нахождения собственных значений матрицы.

Пусть нам необходимо найти характеристический многочлен $D(\lambda) = \det(A - \lambda E)$ квадратной матрицы $A = A_n$ порядка n . Трудность построения такого многочлена

$$D_n(\lambda) = |A_n - \lambda E_n|$$

возрастает с увеличением n . При $n=2$, например, эта задача решается еще совсем просто. Безусловно полезной была бы индуктивная конструкция, посредством которой по характеристическому многочлену $D_{n-1}(\lambda)$

квадратной матрицы A_{n-1} порядка $n-1$ можно было бы построить характеристический многочлен $D_n(\lambda)$ матрицы A_n , полученной окаймлением матрицы A_{n-1} . Такую конструкцию и дает рассматриваемый метод.

Итак, пусть матрицы A_n и A_{n-1} связаны между собой следующим образом:

$$A = A_n = \begin{bmatrix} A_{n-1} & \bar{u}^{(n-1)} \\ \bar{v}^{(n-1)} & a_{nn} \end{bmatrix}. \quad (3.4.2)$$

Здесь

$$\bar{v}^{(n-1)} = (a_{n1}, a_{n2}, \dots, a_{n \ n-1}),$$

$$\bar{u}^{(n-1)} = (a_{1n}, a_{2n}, \dots, a_{n-1 \ n})'.$$

Введем в рассмотрение матрицу $Q(\lambda) = Q_n(\lambda) = (Q_{ij})'$, союзную для матрицы $A_n - \lambda E_n$. По определению союзной матрицы

$$Q_{nn}(\lambda) = D_{n-1}(\lambda).$$

Произведем разбиение на клетки матрицы $Q_n(\lambda)$ аналогично только что выполненному разбиению (3.4.2) матрицы A_n :

$$Q_n(\lambda) = \begin{bmatrix} Q_{n-1}(\lambda) & \bar{g}^{(n-1)}(\lambda) \\ \bar{h}^{(n-1)}(\lambda) & D_{n-1}(\lambda) \end{bmatrix}.$$

Здесь

$$\bar{h}^{(n-1)}(\lambda) = (Q_{1n}(\lambda), Q_{2n}(\lambda), \dots, Q_{n-1 \ n}(\lambda)),$$

$$\bar{g}^{(n-1)}(\lambda) = (Q_{n1}(\lambda), Q_{n2}(\lambda), \dots, Q_{n \ n-1}(\lambda))',$$

а $D_{n-1}(\lambda)$ есть характеристический многочлен матрицы A_{n-1} .

Как известно из алгебры, между матрицей A и ее союзной матрицей C существует следующая связь:

$$AC = |A| E.$$

В нашем случае последнее равенство принимает вид

$$(A_n - \lambda E_n) Q_n(\lambda) = D_n(\lambda) E_n$$

или

$$\begin{bmatrix} A_{n-1} - \lambda E_{n-1} & \bar{u}^{(n-1)} \\ \bar{v}^{(n-1)} & a_{nn} - \lambda \end{bmatrix} \cdot \begin{bmatrix} Q_{n-1}(\lambda) & \bar{g}^{(n-1)}(\lambda) \\ \bar{h}^{(n-1)}(\lambda) & D_{n-1}(\lambda) \end{bmatrix} = D_n(\lambda) E_n.$$

Отсюда, в частности, вытекают следующие равенства:

$$\begin{aligned}(A_{n-1} - \lambda E_{n-1}) \bar{g}^{(n-1)}(\lambda) + \bar{u}^{(n-1)} D_{n-1}(\lambda) &= \bar{0}, \\ \bar{v}^{(n-1)} \bar{g}^{(n-1)}(\lambda) + (a_{nn} - \lambda) D_{n-1}(\lambda) &= D_n(\lambda).\end{aligned}$$

Первое из этих равенств позволяет найти вектор $\bar{g}^{(n-1)}(\lambda)$, после чего второе из них дает возможность построить интересующий нас характеристический многочлен $D_n(\lambda)$ матрицы A_n . При этом обычно первое из упомянутых равенств переписывают в виде

$$\lambda \bar{g}^{(n-1)}(\lambda) = A_{n-1} \bar{g}^{(n-1)}(\lambda) + \bar{u}^{(n-1)} D_{n-1}(\lambda)$$

и вектор $\bar{g}^{(n-1)}(\lambda)$ находят последовательно по слагаемым путем сравнения коэффициентов при одинаковых степенях λ в правой и левой частях последнего векторного равенства. Начинают этот процесс с коэффициента при λ^{n-2} в $\bar{g}^{(n-1)}(\lambda)$, равного, очевидно, вектору $\bar{u}^{(n-1)}$.

Таким образом, исходя из непосредственно вычисляемого характеристического многочлена $D_2(\lambda)$, можно последовательно находить многочлены $D_3(\lambda)$, $D_4(\lambda)$, ..., $D_{n-1}(\lambda)$ и $D_n(\lambda) = D(\lambda)$.

3.4.5. Эскалаторный метод

Этот метод также носит индуктивный характер и представляет собой совокупность правил, посредством которых по известным собственным значениям и собственным векторам матрицы A_{n-1} и ее транспонированной можно построить уравнение, корнями которого будут являться собственные значения матрицы A_n , полученной окаймлением матрицы A_{n-1} , а затем по найденным собственным значениям матрицы A_n найти соответствующие собственные векторы этой матрицы и ее транспонированной. В отличие от рассмотренного выше метода окаймления эскалаторный метод требует на каждом этапе фактического решения соответствующего характеристического уравнения.

Для простоты изложения идеи метода рассмотрим только случай симметричной матрицы, при этом будем предполагать, что все ее собственные значения различны.

Пусть симметричные матрицы $A = A_n$ и A_{n-1} связаны соотношением (3.4.2), при этом вектор-столбец $\bar{u}^{(n-1)}$ получен транспонированием вектора-строки $\bar{v}^{(n-1)}$. Будем считать собственные значения

$$\lambda_1 < \lambda_2 < \dots < \lambda_{n-1}$$

матрицы A_{n-1} и соответствующие им ортонормированные собственные векторы

$$\bar{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_{n-1}^{(i)})' \quad (i = 1, 2, \dots, n-1)$$

известными. Для них справедливо равенство

$$A_{n-1} X_{n-1} = X_{n-1} \Lambda_{n-1}, \quad (3.4.3)$$

где

$$X_{n-1} = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(n-1)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(n-1)} \\ \dots & \dots & \dots & \dots \\ x_{n-1}^{(1)} & x_{n-1}^{(2)} & \dots & x_{n-1}^{(n-1)} \end{bmatrix}, \quad \Lambda_{n-1} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \lambda_{n-1} \end{bmatrix}.$$

Собственный вектор \bar{y} матрицы $A = A_n$ станем искать в виде

$$\bar{y} = \begin{bmatrix} X_{n-1} \bar{z}^{(n-1)} \\ \alpha \end{bmatrix}; \quad (3.4.4)$$

где $\bar{z}^{(n-1)}$ — некоторый $(n-1)$ -мерный вектор-столбец, а α — число.

Учитывая равенства (3.4.2), (3.4.4), из требования

$$A\bar{y} = \lambda\bar{y}$$

получаем:

$$A_{n-1}X_{n-1}\bar{z}^{(n-1)} + \alpha\bar{u}^{(n-1)} = \lambda X_{n-1}\bar{z}^{(n-1)}, \quad (3.4.5)$$

$$\bar{v}^{(n-1)}X_{n-1}\bar{z}^{(n-1)} + \alpha a_{nn} = \lambda\alpha. \quad (3.4.6)$$

Уравнение (3.4.5) с учетом равенства (3.4.3) может быть записано в виде

$$X_{n-1}\Lambda_{n-1}\bar{z}^{(n-1)} + \alpha\bar{u}^{(n-1)} = \lambda X_{n-1}\bar{z}^{(n-1)}.$$

Умножая последнее уравнение на матрицу X'_{n-1} слева и учитывая равенство $X'_{n-1}X_{n-1} = E_{n-1}$ (условие ортонормированности собственных векторов матрицы A_{n-1}), находим

$$\Lambda_{n-1}\bar{z}^{(n-1)} + \alpha X'_{n-1}\bar{u}^{(n-1)} = \lambda\bar{z}^{(n-1)}.$$

Отсюда находим вектор $\bar{z}^{(n-1)}$:

$$\bar{z}^{(n-1)} = \alpha(\lambda E_{n-1} - \Lambda_{n-1})^{-1} X'_{n-1}\bar{u}^{(n-1)}. \quad (3.4.7)$$

Подставим это выражение для $\bar{z}^{(n-1)}$ в равенство (3.4.4):

$$\bar{y} = \alpha \begin{bmatrix} X_{n-1}(\lambda E_{n-1} - \Lambda_{n-1})^{-1} X'_{n-1}\bar{u}^{(n-1)} \\ 1 \end{bmatrix}.$$

Так как собственный вектор матрицы определен с точностью до постоянного множителя, то мы можем выбрать α произвольным, отличным от нуля числом. Тогда последнее равенство дает нам возможность найти собственный вектор матрицы $A = A_n$, принадлежащий ее собственному значению λ , если, конечно, последнее известно. Значение же λ можно найти из уравнения (3.4.6), если подставить туда вместо вектора $\bar{z}^{(n-1)}$ его выражение по формуле (3.4.7). Произведя после такой подстановки сокращение на $\alpha \neq 0$, можем записать:

$$\bar{v}^{(n-1)}X_{n-1}(\lambda E_{n-1} - \Lambda_{n-1})^{-1} X'_{n-1}\bar{u}^{(n-1)} = \lambda - a_{nn}.$$

Так как вектор-столбец $\bar{u}^{(n-1)}$ получен транспонированием вектора-строки $\bar{v}^{(n-1)}$, то в скалярной форме это уравнение принимает следующий вид:

$$\sum_{i=1}^{n-1} \frac{(\bar{v}^{(n-1)}, \bar{x}^{(i)})^2}{\lambda - \lambda_i} = \lambda - a_{nn}.$$

Последнее уравнение обычно называют эскалаторной формой характеристического уравнения матрицы $A = A_n$, полученной окаймлением матрицы A_{n-1} . Найдя корни этого уравнения, мы получим все n собственных значений матрицы A_n . Процесс разыскания корней эскалаторного уравнения матрицы A_n значительно облегчается тем обстоятельством, что они отделены известными собственными значениями матрицы A_{n-1} : имеется точно один корень $\lambda < \lambda_1$, точно один корень между каждой парой λ_i, λ_{i+1} ($i=1, 2, \dots, n-2$) последовательных собственных значений матрицы A_{n-1} и точно один корень $\lambda > \lambda_n$. Действительно, при изменении λ от $-\infty$ до λ_1 левая часть эскалаторного уравнения убывает от 0 до $-\infty$, при изменении λ от λ_1 до λ_2 она убывает от $+\infty$ до $-\infty$ и т. д., при изменении λ от λ_n до $+\infty$ левая часть рассматриваемого уравнения убывает от $+\infty$ до 0. Правая же часть этого уравнения всюду линейно возрастает. Это и доказывает разделение корней эскалаторного уравнения. Такое обстоятельство значительно облегчает задачу вычисления собственных значений матрицы A , которые обычно находятся по методу Ньютона. Эскалаторная форма характеристического уравнения удобна для применения метода Ньютона, так как вычисление значений функций

$$f(\lambda) = \lambda - a_{nn} - \sum_{i=1}^{n-1} \frac{(\bar{v}^{(n-1)}, \bar{x}^{(i)})^2}{\lambda - \lambda_i}$$

и

$$f'(\lambda) = 1 + \sum_{i=1}^{n-1} \frac{(\bar{v}^{(n-1)}, \bar{x}^{(i)})^2}{(\lambda - \lambda_i)^2}$$

нужно производить по очень близким формулам.

3.4.6. Метод ортогонализации

Этот метод, подобно методу Крылова, основан на построении равной нулю линейной комбинации векторов, полученных последовательным итерированием с помощью исходной матрицы A произвольного вектора $\bar{c}^{(0)} \neq \bar{0}$. Однако если в методе Крылова построение такой линейной комбинации связано с решением системы линейных алгебраических уравнений, то в рассматриваемом методе для этого применяется процесс ортогонализации. Как и знакомый уже нам метод ортогонализации решения систем линейных алгебраических уравнений (см. п. 2.5.1), излагаемый здесь метод ортогонализации решения проблемы собственных значений матрицы предполагает предварительное построение системы взаимно ортогональных векторов. Построение таких векторов ведется последова-

По исходному вектору $\bar{c}^{(0)} \neq \bar{0}$ и его итерации $A\bar{c}^{(0)}$ строим вектор

$$\bar{c}^{(1)} = A\bar{c}^{(0)} - g_{10}\bar{c}^{(0)},$$

ортогональный к вектору $\bar{c}^{(0)}$. Это всегда можно сделать. Условие ортогональности $(\bar{c}^{(1)}, \bar{c}^{(0)}) = 0$ позволяет подобрать нужный коэффициент g_{10} :

$$g_{10} = \frac{(A\bar{c}^{(0)}, \bar{c}^{(0)})}{(\bar{c}^{(0)}, \bar{c}^{(0)})}.$$

Если окажется, что $\bar{c}^{(1)} = \bar{0}$, то это будет означать, что векторы $\bar{c}^{(0)}$ и $A\bar{c}^{(0)}$ линейно зависимы. Тогда многочлен

$$\varphi_1(\lambda) = \lambda - g_{10}$$

будет минимальным аннулирующим вектор $\bar{c}^{(0)}$ многочленом матрицы A , так как $\bar{c}^{(0)} \neq \bar{0}$, а

$$\varphi_1(A)\bar{c}^{(0)} = (A - g_{10}E)\bar{c}^{(0)} = \bar{c}^{(1)} = \bar{0}.$$

Корни такого многочлена, как мы знаем, являются собственными значениями матрицы. Если же $\bar{c}^{(1)} \neq \bar{0}$, то строим вектор $A\bar{c}^{(1)}$ и составляем новый вектор

$$\bar{c}^{(2)} = A\bar{c}^{(1)} - g_{21}\bar{c}^{(1)} - g_{20}\bar{c}^{(0)},$$

ортогональный к векторам $\bar{c}^{(1)}$ и $\bar{c}^{(0)}$. Требования ортогональности $(\bar{c}^{(2)}, \bar{c}^{(1)}) = 0$ и $(\bar{c}^{(2)}, \bar{c}^{(0)}) = 0$ дают нужные коэффициенты g_{21} и g_{20} :

$$g_{21} = \frac{(A\bar{c}^{(1)}, \bar{c}^{(1)})}{(\bar{c}^{(1)}, \bar{c}^{(1)})}, \quad g_{20} = \frac{(A\bar{c}^{(0)}, \bar{c}^{(0)})}{(\bar{c}^{(0)}, \bar{c}^{(0)})}.$$

Если построенный вектор $\bar{c}^{(2)}$ будет нулевым, то равенство

$$\bar{0} = A\bar{c}^{(1)} - g_{21}\bar{c}^{(1)} - g_{20}\bar{c}^{(0)}$$

или

$$\bar{0} = [(A - g_{21}E)(A - g_{10}E) - g_{20}E]\bar{c}^{(0)}$$

будет давать нулевую линейную комбинацию

$$A^2\bar{c}^{(0)} - (g_{10} + g_{21})A\bar{c}^{(0)} - (g_{20} - g_{21}g_{10})\bar{c}^{(0)} = \bar{0}$$

векторов $A^2\bar{c}^{(0)}$, $A\bar{c}^{(0)}$, $\bar{c}^{(0)}$. Тогда многочлен

$$\varphi_2(\lambda) = (\lambda - g_{21})(\lambda - g_{10}) - g_{20} = (\lambda - g_{21})\varphi_1(\lambda) - g_{20}$$

будет делителем минимального многочлена матрицы A и его корни, следовательно, будут собственными значениями этой матрицы. Если же окажется, что $\bar{c}^{(2)} \neq \bar{0}$, то процесс ортогонализации следует продолжить.

Пусть выполнено $m-1$ шагов процесса ортогонализации и по исходному вектору $\bar{c}^{(0)} \neq \bar{0}$ построены ненулевые векторы $\bar{c}^{(1)}$, $\bar{c}^{(2)}$, ..., $\bar{c}^{(m-1)}$, удовлетворяющие условиям $(\bar{c}^{(i)}, \bar{c}^{(j)}) = 0$ при $i \neq j$ ($i, j = 0, 1, 2, \dots, m-1$). Тогда составляем вектор

$$\bar{c}^{(m)} = A\bar{c}^{(m-1)} - g_{m-1}\bar{c}^{(m-1)} - g_{m-2}\bar{c}^{(m-2)} - \dots - g_{m0}\bar{c}^{(0)},$$

при этом коэффициенты g_{m-1} , g_{m-2} , ..., g_{m0} подбираем так, чтобы этот вектор был ортогонален к каждому из векторов $\bar{c}^{(0)}$, $\bar{c}^{(1)}$, ..., $\bar{c}^{(m-1)}$. Требования ортогональности $(\bar{c}^{(m)}, \bar{c}^{(i)}) = 0$ ($i = 0, 1, 2, \dots, m-1$) дают нужные нам коэффициенты g_{mi} :

$$g_{mi} = \frac{(A\bar{c}^{(m-1)}, \bar{c}^{(i)})}{(\bar{c}^{(i)}, \bar{c}^{(i)})} \quad (i = 0, 1, 2, \dots, m-1).$$

Одновременно с построением векторов $\bar{c}^{(0)}$, $\bar{c}^{(1)}$, ..., $\bar{c}^{(m)}$ строим многочлены

$$\begin{aligned} \varphi_0(\lambda) &= 1, \\ \varphi_1(\lambda) &= (\lambda - g_{10})\varphi_0(\lambda), \\ \varphi_2(\lambda) &= (\lambda - g_{21})\varphi_1(\lambda) - g_{20}\varphi_0(\lambda), \\ \varphi_3(\lambda) &= (\lambda - g_{32})\varphi_2(\lambda) - g_{31}\varphi_1(\lambda) - g_{30}\varphi_0(\lambda), \\ &\dots \\ \varphi_m(\lambda) &= (\lambda - g_{m-1})\varphi_{m-1}(\lambda) - g_{m-2}\varphi_{m-2}(\lambda) - \dots - g_{m0}\varphi_0(\lambda). \end{aligned}$$

Поскольку в рассматриваемом n -мерном векторном пространстве имеется не более n взаимно ортогональных векторов, то на каком-то m -м ($m \leq n$) шаге процесса ортогонализации обязательно получим нулевой вектор $\bar{c}^{(m)}$. Тогда равенство

$$\bar{0} = A\bar{c}^{(m-1)} - g_{m-1}\bar{c}^{(m-1)} - g_{m-2}\bar{c}^{(m-2)} - \dots - g_{m0}\bar{c}^{(0)}$$

будет означать линейную зависимость векторов

$$\bar{c}^{(0)}, A\bar{c}^{(0)}, A^2\bar{c}^{(0)}, \dots, A^m\bar{c}^{(0)}.$$

Поэтому соответствующий многочлен $\varphi_m(\lambda)$ будет делителем минимального многочлена матрицы A . Если такое обстоятельство встретится лишь при $m=n$, то многочлен $\varphi_n(\lambda)$ будет являться собственным многочленом матрицы A . В случае же $m < n$ мы будем иметь лишь делитель собственного многочлена и сможем, вообще говоря, найти лишь часть собственных значений матрицы. Для отыскания недостающих собственных значений в этом случае приходится брать новый начальный вектор. Его выбирают ортогональным к векторам $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(m-1)}$ и весь процесс повторяют заново.

Система взаимно ортогональных векторов $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(m-1)}$, построенных по методу ортогонализации при разыскании минимального аннулирующего вектор $\bar{c}^{(0)}$ многочлена $\Phi_m(\lambda)$ матрицы A , может быть использована и для нахождения собственных векторов этой матрицы.

Пусть λ_i — корень многочлена $\varphi_m(\lambda)$. Подобно тому, как мы поступаем при нахождении собственных векторов в методе Крылова, будем искать собственный вектор $\bar{x}^{(i)}$ матрицы A , принадлежащий данному собственному значению λ_i , в виде линейной комбинации векторов $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(m-1)}$:

$$\bar{x}^{(i)} = \beta_{i1} \bar{c}^{(m-1)} + \beta_{i2} \bar{c}^{(m-2)} + \dots + \beta_{im} \bar{c}^{(0)}.$$

Коэффициенты β_{ij} ($j=1, 2, \dots, m$) подберем так, чтобы удовлетворить условию

$$A\bar{x}^{(i)} = \lambda\bar{x}^{(i)}.$$

Умножая записанную линейную комбинацию на матрицу A и учитывая равенства

$$A\bar{c}^{(j)} = \bar{c}^{(j+1)} + g_{j+1} \bar{c}^{(j)} + g_{j+1} \bar{c}^{(j-1)} + \dots + g_{j+1} \bar{c}^{(0)} \quad (j=0, 1, 2, \dots, m-1)$$

и требование

$$A\bar{x}^{(i)} = \lambda_i \bar{x}^{(i)},$$

ПОЛУЧИМ:

[illegible]

или

[illegible]

Таким образом, собственный вектор $\bar{x}^{(i)}$ матрицы A , принадлежащий собственному значению λ_i , может быть записан в следующем виде:

$$\bar{x}^{(i)} = \psi_0(\lambda_i) \bar{c}^{(m-1)} + \psi_1(\lambda_i) \bar{c}^{(m-2)} + \dots + \psi_{m-1}(\lambda_i) \bar{c}^{(0)}. \quad (3.4.9)$$

Изложенный выше подход к нахождению собственных векторов матрицы в методе ортогонализации близок рассмотренному ранее (п. 3.2.3) способу нахождения собственных векторов матрицы в методе Крылова и применим при любом $m \leq n$. В случае же $m = n$ можно предложить, например, и несколько иной путь к решению этой задачи, более близкий по идее к подходу, характерному для регулярного случая метода Данилевского (см. п. 3.3.2).

Непосредственно проверяется, что система векторных равенств

$$\begin{aligned} \bar{c}^{(1)} &= A\bar{c}^{(0)} - g_{10}\bar{c}^{(0)}, \\ \bar{c}^{(2)} &= A\bar{c}^{(1)} - g_{21}\bar{c}^{(1)} - g_{20}\bar{c}^{(0)}, \\ &\dots \\ \bar{c}^{(n-1)} &= A\bar{c}^{(n-2)} - g_{n-1, n-2}\bar{c}^{(n-2)} - g_{n-1, n-3}\bar{c}^{(n-3)} - \dots - g_{n-1, 0}\bar{c}^{(0)}, \\ \bar{0} = \bar{c}^{(n)} &= A\bar{c}^{(n-1)} - g_{n, n-1}\bar{c}^{(n-1)} - g_{n, n-2}\bar{c}^{(n-2)} - \dots - g_{n, 0}\bar{c}^{(0)} \end{aligned}$$

равносильна матричному равенству

$$AC - CG = 0, \quad (3.4.10)$$

где

$$C = [\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(n-1)}],$$

$$G = \begin{bmatrix} g_{10} & g_{20} & g_{30} & \dots & g_{n-1, 0} & g_{n, 0} \\ 1 & g_{21} & g_{31} & \dots & g_{n-1, 1} & g_{n, 1} \\ 0 & 1 & g_{32} & \dots & g_{n-1, 2} & g_{n, 2} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 & g_{n, n-1} \end{bmatrix}.$$

Так как векторы $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(n-1)}$ линейно независимы, то существует матрица C^{-1} и последнее матричное равенство может быть приведено к виду

$$A = CGC^{-1}.$$

Значит, матрица A подобна матрице G и собственный вектор $\bar{x}^{(i)}$ матрицы A , принадлежащий собственному значению λ_i , связан соотношением

$$\bar{x}^{(i)} = C\bar{y}^{(i)}$$

с собственным вектором $\bar{y}^{(i)} = (y_1^{(i)}, y_2^{(i)}, \dots, y_n^{(i)})'$ матрицы G , соответствующим тому же собственному значению (см., например, п. 3.3.2). Собственные же векторы матрицы G находятся просто. Действительно, записав равенство

$$G\bar{y}^{(i)} = \lambda_i \bar{y}^{(i)}$$

покоординатно:

непосредственно следует, что $g_{i+1\ j}=0$, если $j < i-1$. Поэтому для симметричных матриц равенства

$$\begin{aligned}\bar{c}^{(i+1)} &= A\bar{c}^{(i)} - g_{i+1\ i}\bar{c}^{(i)} - g_{i+1\ i-1}\bar{c}^{(i-1)} - \dots - g_{i+1\ 0}\bar{c}^{(0)} \\ (i &= 0, 1, 2, \dots, m-1; \ m \leq n)\end{aligned}$$

принимают вид

$$\bar{c}^{(i+1)} = A\bar{c}^{(i)} - g_{i+1\ i}\bar{c}^{(i)} - g_{i+1\ i-1}\bar{c}^{(i-1)}$$

и вид многочленов

$$\begin{aligned}\varphi_{i+1}(\lambda) &= (\lambda - g_{i+1\ i})\varphi_i(\lambda) - g_{i+1\ i-1}\varphi_{i-1}(\lambda) - \dots - g_{i+1\ 0}\varphi_0(\lambda) \\ (i &= 0, 1, 2, \dots, m-1; \ m \leq n)\end{aligned}$$

также упрощается:

$$\varphi_{i+1}(\lambda) = (\lambda - g_{i+1\ i})\varphi_i(\lambda) - g_{i+1\ i-1}\varphi_{i-1}(\lambda).$$

Это позволяет значительно упростить и вычислительную схему метода. Поэтому метод ортогонализации в применении к симметричным матрицам называют обычно методом минимальных итераций.

Подобного же упрощения можно добиться и для случая несимметричной матрицы, заменив процесс ортогонализации процессом биортогонализации, который мы сейчас и рассмотрим.

Изберем два начальных вектора $\bar{c}^{(0)}$ и $\bar{b}^{(0)}$ и по ним построим векторы

$$\bar{c}^{(1)} = A\bar{c}^{(0)} - g_{10}\bar{c}^{(0)}$$

и

$$\bar{b}^{(1)} = A'\bar{b}^{(0)} - h_{10}\bar{b}^{(0)}$$

такие, чтобы выполнялись условия

$$(\bar{c}^{(1)}, \bar{b}^{(0)}) = (\bar{b}^{(1)}, \bar{c}^{(0)}) = 0.$$

Если исходные векторы $\bar{c}^{(0)}$ и $\bar{b}^{(0)}$ не были ортогональны, то искомые векторы $\bar{c}^{(1)}$ и $\bar{b}^{(1)}$ всегда можно построить, при этом

$$g_{10} = \frac{(A\bar{c}^{(0)}, \bar{b}^{(0)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = \frac{(\bar{c}^{(0)}, A'\bar{b}^{(0)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = h_{10}.$$

Будем предполагать условие $(\bar{c}^{(0)}, \bar{b}^{(0)}) \neq 0$ выполненным и по найденным векторам $\bar{c}^{(1)}$ и $\bar{b}^{(1)}$ построим новые векторы

$$\bar{c}^{(2)} = A\bar{c}^{(1)} - g_{21}\bar{c}^{(1)} - g_{20}\bar{c}^{(0)}$$

и

$$\bar{b}^{(2)} = A'\bar{b}^{(1)} - h_{21}\bar{b}^{(1)} - h_{20}\bar{b}^{(0)}$$

такие, что

$$(\bar{c}^{(2)}, \bar{b}^{(1)}) = (\bar{c}^{(2)}, \bar{b}^{(0)}) = (\bar{b}^{(2)}, \bar{c}^{(1)}) = (\bar{b}^{(2)}, \bar{c}^{(0)}) = 0.$$

Такое построение возможно, если векторы $\bar{c}^{(1)}$ и $\bar{b}^{(1)}$ также неортогональны, при этом получим:

$$\begin{aligned} g_{21} &= \frac{(A\bar{c}^{(1)}, \bar{b}^{(1)})}{(\bar{c}^{(1)}, \bar{b}^{(1)})} = \frac{(\bar{c}^{(1)}, A'\bar{b}^{(1)})}{(\bar{c}^{(1)}, \bar{b}^{(1)})} = h_{21}, \\ g_{20} &= \frac{(A\bar{c}^{(1)}, \bar{b}^{(0)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = \frac{(\bar{c}^{(1)}, A'\bar{b}^{(0)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = \\ &= \frac{(\bar{c}^{(1)}, \bar{b}^{(1)} + h_{10}\bar{b}^{(0)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = \frac{(\bar{c}^{(1)}, \bar{b}^{(1)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = \frac{(A\bar{c}^{(0)} - g_{10}\bar{c}^{(0)}, \bar{b}^{(1)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = \\ &= \frac{(A\bar{c}^{(0)}, \bar{b}^{(1)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = \frac{(\bar{c}^{(0)}, A'\bar{b}^{(1)})}{(\bar{c}^{(0)}, \bar{b}^{(0)})} = h_{20}. \end{aligned}$$

Предположим, что

$$(\bar{c}^{(0)}, \bar{b}^{(0)}) \neq 0, \quad (\bar{c}^{(1)}, \bar{b}^{(1)}) \neq 0, \quad (\bar{c}^{(2)}, \bar{b}^{(2)}) \neq 0,$$

и продолжим процесс биортогонализации. Пусть подобным же образом построены векторы

$$\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(i)},$$

$$\bar{b}^{(0)}, \bar{b}^{(1)}, \dots, \bar{b}^{(i)},$$

такие, что

$$(\bar{b}^{(j)}, \bar{c}^{(k)}) = 0 \quad (j \neq k)$$

и

$$(\bar{b}^{(j)}, \bar{c}^{(j)}) \neq 0 \quad (j, k = 0, 1, 2, \dots, i).$$

Тогда следующая пара векторов $\bar{c}^{(i+1)}, \bar{b}^{(i+1)}$ находится по формулам

$$\left. \begin{aligned} \bar{c}^{(i+1)} &= A\bar{c}^{(i)} - g_{i+1 \ i} \bar{c}^{(i)} - g_{i+1 \ i-1} \bar{c}^{(i-1)} - \dots - g_{i+1 \ 0} \bar{c}^{(0)}, \\ \bar{b}^{(i+1)} &= A'\bar{b}^{(i)} - h_{i+1 \ i} \bar{b}^{(i)} - h_{i+1 \ i-1} \bar{b}^{(i-1)} - \dots - h_{i+1 \ 0} \bar{b}^{(0)}, \end{aligned} \right\} \quad (3.4.11)$$

при этом коэффициенты $g_{i+1 \ j}$ и $h_{i+1 \ j}$ ($j=0, 1, 2, \dots, i$) выбираются такими, чтобы выполнялись условия

$$(\bar{c}^{(i+1)}, \bar{b}^{(j)}) = (\bar{b}^{(i+1)}, \bar{c}^{(j)}) = 0 \quad (j=0, 1, 2, \dots, i).$$

Удовлетворяя эти требования, находим, что

$$\begin{aligned} g_{i+1 \ j} &= \frac{(A\bar{c}^{(i)}, \bar{b}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \frac{(\bar{c}^{(i)}, A'\bar{b}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \frac{(\bar{c}^{(i)}, \bar{b}^{(j+1)} + h_{j+1 \ j} \bar{b}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \\ &= \frac{(\bar{c}^{(i)}, \bar{b}^{(j+1)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} + h_{j+1 \ j} \frac{(\bar{c}^{(i)}, \bar{b}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \begin{cases} h_{i+1 \ i} & \text{при } j=i, \\ \frac{(\bar{c}^{(i)}, \bar{b}^{(i)})}{(\bar{c}^{(i-1)}, \bar{b}^{(i-1)})} & \text{при } j=i-1, \\ 0 & \text{при } j < i-1, \end{cases} \\ h_{i+1 \ j} &= \frac{(A'\bar{b}^{(i)}, \bar{c}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \frac{(\bar{b}^{(i)}, A\bar{c}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \frac{(\bar{b}^{(i)}, \bar{c}^{(j+1)} + g_{j+1 \ j} \bar{c}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \\ &= \frac{(\bar{b}^{(i)}, \bar{c}^{(j+1)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} + g_{j+1 \ j} \frac{(\bar{b}^{(i)}, \bar{c}^{(j)})}{(\bar{c}^{(j)}, \bar{b}^{(j)})} = \\ &= \begin{cases} g_{i+1 \ i} & \text{при } j=i, \\ \frac{(\bar{c}^{(i)}, \bar{b}^{(i)})}{(\bar{c}^{(i-1)}, \bar{b}^{(i-1)})} = g_{i+1 \ i-1} & \text{при } j=i-1, \\ 0 & \text{при } j < i-1. \end{cases} \end{aligned}$$

Таким образом, равенства (3.4.11), посредством которых строятся биортогональные системы векторов, значительно упрощаются и принимают следующий вид:

$$\begin{aligned} \bar{c}^{(i+1)} &= A\bar{c}^{(i)} - g_{i+1 \ i} \bar{c}^{(i)} - g_{i+1 \ i-1} \bar{c}^{(i-1)}, \\ \bar{b}^{(i+1)} &= A'\bar{b}^{(i)} - g_{i+1 \ i} \bar{b}^{(i)} - g_{i+1 \ i-1} \bar{b}^{(i-1)}. \end{aligned}$$

Возьмем произвольный вектор $\bar{c}^{(0)} \neq \bar{0}$ и по нему построим вектор

$$\bar{c}^{(1)} = A\bar{c}^{(0)} - g_{10}\bar{c}^{(0)},$$

при этом коэффициент g_{10} подберем так, чтобы первая координата вектора $\bar{c}^{(1)}$ была равна нулю. По векторам $\bar{c}^{(0)}$ и $\bar{c}^{(1)}$ строим, далее, вектор

$$\bar{c}^{(2)} = A\bar{c}^{(1)} - g_{21}\bar{c}^{(1)} - g_{20}\bar{c}^{(0)},$$

выбирая коэффициенты g_{21} и g_{20} из условия равенства нулю двух первых координат вектора $\bar{c}^{(2)}$. Подобным же образом продолжаем процесс построения векторов $\bar{c}^{(i)}$, первые i координат которых будут нулевыми. При этом каждый последующий вектор получается итерированием матрицей A предыдущего вектора с последующей коррекцией результата посредством добавления подходящей линейной комбинации всех предшествующих векторов, т. е.

$$\begin{aligned} \bar{c}^{(i+1)} &= (0, 0, \dots, 0, c_{i+2}^{(i+1)}, c_{i+3}^{(i+1)}, \dots, c_n^{(i+1)})' = \\ &= A\bar{c}^{(i)} - g_{i+1, i}\bar{c}^{(i)} - g_{i+1, i-1}\bar{c}^{(i-1)} - \dots - g_{i+1, 0}\bar{c}^{(0)}. \end{aligned}$$

Такой процесс построения по вектору $\bar{c}^{(0)}$ векторов $\bar{c}^{(1)}, \bar{c}^{(2)}, \dots, \bar{c}^{(n-1)}$ не всегда удается осуществить до конца. Если окажется, что у вектора $\bar{c}^{(i+1)}$ равна нулю координата $c_{i+2}^{(i+1)}$, то естественное течение процесса нарушается. Будем пока иметь в виду лишь регулярный случай, т. е. будем предполагать, что

$$c_1^{(0)} \neq 0, c_2^{(1)} \neq 0, \dots, c_n^{(n-1)} \neq 0.$$

Тогда векторы $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(n-1)}$ будут, очевидно, линейно независимы и матрица

$$C = [\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(n-1)}]$$

будет неособенной. При этом вектор $\bar{c}^{(i)}$ ($i=0, 1, 2, \dots, n$) может быть представлен в виде

$$\bar{c}^{(i)} = \varphi_i(A)\bar{c}^{(0)},$$

где, как и в случае метода ортогонализации, многочлены $\varphi_i(\lambda)$ строятся последовательно по формулам

$$\left. \begin{aligned} \varphi_0(\lambda) &= 1, \\ \varphi_1(\lambda) &= (\lambda - g_{10})\varphi_0(\lambda), \\ \varphi_2(\lambda) &= (\lambda - g_{21})\varphi_1(\lambda) - g_{20}\varphi_0(\lambda), \\ \varphi_3(\lambda) &= (\lambda - g_{32})\varphi_2(\lambda) - g_{31}\varphi_1(\lambda) - g_{30}\varphi_0(\lambda), \\ &\dots \\ \varphi_i(\lambda) &= (\lambda - g_{i \ i-1})\varphi_{i-1}(\lambda) - g_{i \ i-2}\varphi_{i-2}(\lambda) - \dots - g_{i0}\varphi_0(\lambda). \end{aligned} \right\} \quad (3.4.12)$$

Так как векторы $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(n-1)}$ линейно независимы, то равенство

$$\varphi_i(A)\bar{c}^{(0)} = \bar{0}$$

невозможно при $i < n$. Полином же $\varphi_n(\lambda)$ по построению удовлетворяет требованию

$$\varphi_n(A)\bar{c}^{(0)} = \bar{c}^{(n)} = \bar{0}.$$

Так как, кроме того, степень этого многочлена равна n и его старший коэффициент равен единице, то он совпадает с собственным многочленом матрицы A .

Таким образом, проведя в регулярном случае n шагов метода Хессенберга, мы сможем построить собственный многочлен данной квадратной матрицы порядка n .

Остановимся несколько подробнее на организации вычислений при нахождении коэффициентов g_{ij} ($j < i, i = 1, 2, \dots, n$).

Как и в случае метода ортогонализации, система векторных равенств

$$\begin{aligned} \bar{c}^{(1)} &= A\bar{c}^{(0)} - g_{10}\bar{c}^{(0)}, \\ \bar{c}^{(2)} &= A\bar{c}^{(1)} - g_{21}\bar{c}^{(1)} - g_{20}\bar{c}^{(0)}, \\ &\dots \\ \bar{c}^{(n-1)} &= A\bar{c}^{(n-2)} - g_{n-1 \ n-2}\bar{c}^{(n-2)} - g_{n-1 \ n-3}\bar{c}^{(n-3)} - \dots - g_{n-1 \ 0}\bar{c}^{(0)}, \\ \bar{0} = \bar{c}^{(n)} &= A\bar{c}^{(n-1)} - g_{n \ n-1}\bar{c}^{(n-1)} - g_{n \ n-2}\bar{c}^{(n-2)} - \dots - g_{n0}\bar{c}^{(0)} \end{aligned}$$

равносильна матричному равенству (3.4.10).

Равенство $AC - CG = 0$ позволяет последовательно находить коэффициенты g_{ij} ($j < i, i = 1, 2, \dots, n$) составляемых линейных комбинаций и координаты $c_i^{(j)}$ ($i = 1, 2, \dots, n; j = 0, 1, 2, \dots, n-1$) векторов $\bar{c}^{(j)}$, при этом для удобства вычислений его обычно представляют в форме

$$(A|C) \left(\frac{C}{-G} \right) = 0,$$

где прямоугольные матрицы $(A|C)$ и $\left(\frac{C}{-G} \right)$ имеют следующий вид:

$$(A|C) = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} & c_1^{(0)} & 0 & \dots & 0 \\ a_{21} & a_{22} & \dots & a_{2n} & c_2^{(0)} & c_2^{(1)} & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} & c_n^{(0)} & c_n^{(1)} & \dots & c_n^{(n-1)} \end{bmatrix},$$

$$\left(\frac{C}{-G} \right) = \begin{bmatrix} c_1^{(0)} & 0 & \dots & 0 \\ c_2^{(0)} & c_2^{(1)} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ c_n^{(0)} & c_n^{(1)} & \dots & c_n^{(n-1)} \\ -g_{10} & -g_{20} & \dots & -g_{n0} \\ -1 & -g_{21} & \dots & -g_{n1} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & -g_{nn-1} \end{bmatrix}.$$

До начала вычислений мы знаем лишь матрицу A и первый столбец матрицы C . Умножая первую строку матрицы $(A|C)$ на первый столбец матрицы $\left(\frac{C}{-G} \right)$ и приравнявая результат умножения нулю, мы получим линейное уравнение для нахождения коэффициента g_{10} . Точно так же умножение остальных строк матрицы $(A|C)$ на первый столбец матрицы $\left(\frac{C}{-G} \right)$ позволяет последовательно определить координаты $c_2^{(1)}, c_3^{(1)}, \dots, c_n^{(1)}$ вектора $\bar{c}^{(1)}$. После этого умножением матрицы $(A|C)$ на второй столбец матрицы $\left(\frac{C}{-G} \right)$ мы последовательно найдем элементы $g_{21}, g_{20}, c_3^{(2)}, c_4^{(2)}, \dots, c_n^{(2)}$. Далее производим последовательное умножение матрицы $(A|C)$ на остальные столбцы матрицы $\left(\frac{C}{-G} \right)$.

Вычисления по такой схеме условно обозначим следующим образом:

$$\left(\begin{array}{c|c} A & C \\ \hline & -G \end{array} \right).$$

Иногда вычисления организуют также по схеме

$$\begin{pmatrix} A' & C \\ C' & -G \end{pmatrix},$$

заменяв умножение строк на столбцы умножением столбцов. Такая схема, несмотря на двойную запись элементов матрицы C , оправдывает себя простотой действий при работе на настольных машинах.

Найдя значения коэффициентов g_{ij} ($j < i, i = 1, 2, \dots, n$), мы сможем по формулам вида (3.4.12) последовательно записать собственный многочлен $\varphi_n(\lambda)$ исходной матрицы A . Этот многочлен можно найти и иначе, если учесть, что (см. равенства (3.4.10)) матрица

$$A = CGC^{-1}$$

подобна матрице G . Так как собственные многочлены подобных матриц совпадают, то можно искать собственный многочлен матрицы A по матрице G . Учитывая специальный вид этой матрицы, для нахождения ее собственного многочлена с успехом можно воспользоваться методом Крылова.

Отметим, кстати, что равенство $A = CGC^{-1}$ позволяет также истолковывать метод Хессенберга как метод приведения данной матрицы A преобразованием подобия $C^{-1}AC$ с треугольной матрицей C к матрице специального вида

$$G = \begin{bmatrix} g_{10} & g_{20} & g_{30} & \dots & g_{n0} \\ 1 & g_{21} & g_{31} & \dots & g_{n1} \\ 0 & 1 & g_{32} & \dots & g_{n2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & g_{nn-1} \end{bmatrix},$$

собственный многочлен и собственные векторы которой находятся сравнительно просто. Это обстоятельство сближает метод Хессенберга с методом Данилевского. В частности, в регулярном случае метода Хессенберга собственные векторы матрицы A можно находить, как и в методе Данилевского, опираясь на известную связь соответствующих собственных векторов подобных матриц (см., например, п. 3.3.2). Как мы уже отмечали и ранее (см. п. 3.4.6), после вычисления собственных значений матрицы G соответствующие собственные векторы этой матрицы находятся просто из условия

$$G\bar{y} = \lambda\bar{y}.$$

Остановимся теперь вкратце на рассмотрении исключительных случаев, которые могут встретиться при реализации алгоритма Хессенберга.

Предположим, что после выполнения i -го шага процесса оказалась равной нулю (кроме первых i координат) также и $(i+1)$ -я координата вектора $\bar{c}^{(i)}$. Тогда естественное течение процесса нарушается.

Если при этом будут нулевыми и все остальные координаты этого вектора, то процесс построения векторов $\bar{c}^{(j)}$ по данному начальному вектору $\bar{c}^{(0)}$ заканчивается, так как уже получена искомая нулевая линейная комбинация векторов

$$\bar{c}^{(0)}, A\bar{c}^{(0)}, A^2\bar{c}^{(0)}, \dots, A^i\bar{c}^{(0)}.$$

Многочлен $\varphi_i(\lambda)$ при этом обладает свойством

$$\varphi_i(A)\bar{c}^{(0)} = \bar{0}$$

и будет, следовательно, делителем собственного многочлена матрицы A . Его корни дадут нам, вообще говоря, лишь часть собственных значений этой матрицы. Разыскание же остальных собственных значений ее связано, как и в случае других методов подобного типа, с выбором нового начального вектора.

Если же хотя бы одна из координат номера $j > i+1$ вектора $\bar{c}^{(i)}$ отлична от нуля, то это является свидетельством неудачного выбора вектора $\bar{c}^{(0)}$. При этом, правда, не обязательно нужно изменять начальный вектор. Процесс можно продолжить. Но, заполняя в матрице C столбец для вектора $\bar{c}^{(i+1)}$, следует поставить нули только на тех местах, на которых за счет добавления линейной комбинации предшествующих столбцов можно фактически добиться нулевых значений. Если в результате дальнейшего продолжения процесса при этом мы получим n ненулевых векторов $\bar{c}^{(0)}, \bar{c}^{(1)}, \dots, \bar{c}^{(n-1)}$, то они, очевидно, также будут линейно независимы и посредством матрицы C , составленной по ним, можно будет осуществить преобразование подобия матрицы A в новую матрицу G , проблема собственных значений которой решается сравнительно просто. Матрица C в этом случае уже не будет треугольной, но будет получаться из треугольной перестановкой столбцов.

3.4.8. Метод Самуэльсона

Укажем еще на один метод нахождения собственного многочлена матрицы, очень близкий по идее методу Крылова.

Вычислительная схема метода Самуэльсона была первоначально получена автором посредством специального преобразования системы линейных дифференциальных уравнений первого порядка с постоянными коэффициентами, связанной с данной матрицей A , к одному дифференциальному уравнению порядка n . Мы рассмотрим здесь лишь краткое алгебраическое обоснование этой схемы.

Изберем произвольный ненулевой вектор

$$\bar{c}^{(0)} = (c_1^{(0)}, c_2^{(0)}, \dots, c_n^{(0)})' = \begin{bmatrix} c_1^{(0)} \\ -1 \\ g^{(0)} \end{bmatrix}$$

и будем итерировать его посредством исходной матрицы A :

$$\begin{aligned} A\bar{c}^{(0)} = \bar{c}^{(1)} &= (c_1^{(1)}, c_2^{(1)}, \dots, c_n^{(1)})' = \begin{bmatrix} c_1^{(1)} \\ -1 \\ g^{(1)} \end{bmatrix}, \quad A^2\bar{c}^{(0)} = \bar{c}^{(2)} = (c_1^{(2)}, c_2^{(2)}, \dots, c_n^{(2)})' = \\ &= \begin{bmatrix} c_1^{(2)} \\ -1 \\ g^{(2)} \end{bmatrix}, \dots, A^n\bar{c}^{(0)} = \bar{c}^{(n)} = (c_1^{(n)}, c_2^{(n)}, \dots, c_n^{(n)})' = \begin{bmatrix} c_1^{(n)} \\ -1 \\ g^{(n)} \end{bmatrix}. \end{aligned}$$

Матрицу $A = A_n$ путем разбиения на клетки представим в следующем виде:

$$A = A_n = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} a_{11} & \bar{u}^{(n-1)} \\ \bar{v}^{(n-1)} & B_{n-1} \end{bmatrix},$$

где

$$\bar{u}^{(n-1)} = (a_{12}, a_{13}, \dots, a_{1n}), \quad \bar{v}^{(n-1)} = (a_{21}, a_{31}, \dots, a_{n1})'.$$

Тогда результат j -й ($j=1, 2, \dots, n$) итерации вектора $\bar{c}^{(0)}$ матрицей A может быть описан посредством следующих равенств:

$$\begin{aligned} c_1^{(j)} &= a_{11}c_1^{(j-1)} + \bar{u}^{(n-1)}\bar{g}^{(j-1)}, \\ \bar{g}^{(j)} &= \bar{v}^{(n-1)}c_1^{(j-1)} + B_{n-1}\bar{g}^{(j-1)}. \end{aligned}$$

Отсюда последовательно находим

$$\begin{aligned} c_1^{(j)} &= a_{11}c_1^{(j-1)} + \bar{u}^{(n-1)}\bar{g}^{(j-1)} = a_{11}c_1^{(j-1)} + \bar{u}^{(n-1)}\bar{v}^{(n-1)}c_1^{(j-2)} + \bar{u}^{(n-1)}B_{n-1}\bar{g}^{(j-2)} = \\ &= a_{11}c_1^{(j-1)} + \bar{u}^{(n-1)}\bar{v}^{(n-1)}c_1^{(j-2)} + \bar{u}^{(n-1)}B_{n-1}\bar{v}^{(n-1)}c_1^{(j-3)} + \bar{u}^{(n-1)}B_{n-1}^2\bar{g}^{(j-3)} = \\ &= \dots = a_{11}c_1^{(j-1)} + \bar{u}^{(n-1)}\bar{v}^{(n-1)}c_1^{(j-2)} + \bar{u}^{(n-1)}B_{n-1}\bar{v}^{(n-1)}c_1^{(j-3)} + \dots + \\ &\quad + \bar{u}^{(n-1)}B_{n-1}^{j-2}\bar{v}^{(n-1)}c_1^{(0)} + \bar{u}^{(n-1)}B_{n-1}^{j-1}\bar{g}^{(0)} \end{aligned}$$

или

$$\begin{aligned} \bar{u}^{(n-1)}B_{n-1}^{j-1}\bar{g}^{(0)} &= c_1^{(j)} - a_{11}c_1^{(j-1)} - \bar{u}^{(n-1)}\bar{v}^{(n-1)}c_1^{(j-2)} - \bar{u}^{(n-1)}B_{n-1}\bar{v}^{(n-1)}c_1^{(j-3)} - \\ &- \dots - \bar{u}^{(n-1)}B_{n-1}^{j-2}\bar{v}^{(n-1)}c_1^{(0)} \quad (j=1, 2, \dots, n). \end{aligned} \quad (3.4.13)$$

Таким образом, мы получаем систему n линейных соотношений между $\bar{g}^{(0)}, c_1^{(n)}, c_1^{(n-1)}, \dots, c_1^{(0)}$, которую можно охарактеризовать матрицей

$$\left[\begin{array}{c|cccccc} \bar{u}^{(n-1)} & 0 & 0 & 0 & \dots & 1 & -a_{11} \\ \bar{u}^{(n-1)}B_{n-1} & 0 & 0 & 0 & & -a_{11} & -\bar{u}^{(n-1)}\bar{v}^{(n-1)} \\ \bar{u}^{(n-1)}B_{n-1}^2 & 0 & 0 & 0 & \dots & -\bar{u}^{(n-1)}\bar{v}^{(n-1)} & -\bar{u}^{(n-1)}B_{n-1}\bar{v}^{(n-1)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \bar{u}^{(n-1)}B_{n-1}^{n-1} & 1 & -a_{11} & -\bar{u}^{(n-1)}\bar{v}^{(n-1)} & & & -\bar{u}^{(n-1)}B_{n-1}^{n-2}\bar{v}^{(n-1)} \end{array} \right]. \quad (3.4.14)$$

Исключив из этих n равенств координаты вектора $\bar{g}^{(0)}$, мы найдем одно линейное соотношение между числами $c_1^{(n)}, c_1^{(n-1)}, \dots, c_1^{(0)}$. Коэффициенты этого соотношения будут постоянными числами, не зависящими от выбора начального вектора $\bar{c}^{(0)}$.

Для собственного многочлена $P(\lambda)$ по теореме Гамильтона — Кели справедливо равенство $P(A) = O$.

Умножив это матричное равенство на вектор $\bar{c}^{(0)}$, можно, в частности, снова получить линейную зависимость между числами $c_1^{(n)}, c_1^{(n-1)}, \dots, c_1^{(0)}$ с постоянными коэффициентами $1, -p_1, -p_2, \dots, -p_n$, не зависящими от выбора начального вектора $\bar{c}^{(0)}$:

$$c_1^{(n)} - p_1 c_1^{(n-1)} - p_2 c_1^{(n-2)} - \dots - p_n c_1^{(0)} = 0$$

(мы выписали здесь лишь связь между первыми координатами векторного равенства $P(A)\bar{c}^{(0)} = \bar{0}$).

Приведенное соотношение между числами $c_1^{(n)}, c_1^{(n-1)}, \dots, c_1^{(0)}$ будет совпадать с соответствующим соотношением, получаемым в процессе исключения из равенств (3.4.13) координат вектора $\bar{g}^{(0)}$, если матрица A такова, что числа $c_1^{(0)}, c_1^{(1)}, \dots, c_1^{(n-1)}$ можно считать независимыми переменными, т. е. если им можно независимо друг от друга придавать произвольные значения, подбирая подходящим образом вектор $\bar{g}^{(0)}$, получаемый из вектора $\bar{c}^{(0)}$ усечением первой координаты. В этом случае сравнение соответствующих коэффициентов этих соотношений даст нам искомые коэффициенты собственного многочлена матрицы A .

Таким образом, вычислительная схема метода Самуэльсона предполагает построение прямоугольной матрицы (3.4.14) и исключение из нее посредством элементарных преобразований вектора-строки $\bar{u}^{(n-1)}B_{n-1}^{n-1}$. Тогда остальные элементы последней строки матрицы (3.4.14) будут давать, вообще говоря, коэффициенты собственного многочлена исходной матрицы A .

§ 3.5. ИТЕРАЦИОННЫЕ МЕТОДЫ НАХОЖДЕНИЯ СОБСТВЕННЫХ ЗНАЧЕНИЙ И СОБСТВЕННЫХ ВЕКТОРОВ МАТРИЦЫ

Этим параграфом мы начинаем рассмотрение так называемой частичной проблемы собственных значений, состоящей в определении обычно одного или нескольких наибольших по модулю собственных значений матрицы и принадлежащих им собственных векторов. Для решения такой проблемы разработано большое число методов, в основу которых положены идеи использования тех или других частных свойств собственных

значений и собственных векторов, например распределение по модулю и кратность собственных значений, ортогональность собственных векторов и т. д. Все методы, предназначенные для решения частичной проблемы, являются итерационными.

3.5.1. Степенной метод для вычисления наибольшего по модулю собственного значения матрицы и соответствующего собственного вектора

Этот метод позволяет находить наибольшее по модулю собственное значение матрицы и принадлежащий ему собственный вектор при помощи вычисления последовательности итераций произвольного вектора матрицей A до тех пор, пока в этой последовательности станет преобладающей одна составляющая в разложении упомянутого вектора по собственным векторам матрицы. Качество указанного итерационного процесса существенно зависит от того, как входит наибольшее по модулю собственное значение матрицы в ее каноническую форму Жордана. Как мы увидим, процесс может усложниться и не привести к цели, если наибольшему по модулю собственному значению матрицы A будут соответствовать нелинейные элементарные делители высокой степени. Поэтому мы остановимся на некоторых простых случаях степенного метода, когда вычисления, как правило, приводят к цели. В частности, мы будем предполагать, что элементарные делители матрицы A линейны. Такое предположение наверное будет выполняться в двух важных частных случаях: 1) матрица A — симметрическая, 2) собственные числа матрицы A различны.

Отдельно будет рассмотрен случай, когда элементарный делитель матрицы A , отвечающий наибольшему по модулю собственному значению, имеет вторую степень.

Итак, пусть A — вещественная матрица, все собственные значения которой имеют линейные элементарные делители. Обозначим собственные значения и соответствующие им собственные векторы матрицы A через

$$\lambda_1, \lambda_2, \dots, \lambda_n$$

и

$$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n.$$

Для определенности записи условимся, что собственные значения матрицы A перенумерованы в порядке невозрастания их модулей, т. е.

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|.$$

Будем рассматривать случай, когда λ_1 — наибольшее по модулю собственное значение, вещественное и простое, и имеют место неравенства

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|. \quad (3.5.1)$$

Выберем произвольный ненулевой вектор $\bar{y}^{(0)}$ и с помощью матрицы A построим итерационную последовательность векторов

$$\bar{y}^{(1)}, \bar{y}^{(2)}, \dots, \bar{y}^{(k)}, \dots$$

по следующему степенному правилу:

$$\bar{y}^{(k)} = A\bar{y}^{(k-1)} = \dots = A^k \bar{y}^{(0)} \quad (k=1, 2, \dots).$$

При сделанных предположениях относительно матрицы A ее собственные векторы образуют полную систему и, следовательно, мы можем записать:

$$\bar{y}^{(0)} = \alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2 + \dots + \alpha_n \bar{x}_n. \quad (3.5.2)$$

Здесь α_i — некоторые числа, среди которых могут быть и равные нулю. Предположим, что $\alpha_1 \neq 0$. Если $\alpha_1 = 0$ и это условие будет каким-либо образом выявлено, то выбор начального вектора $\bar{y}^{(0)}$ следует изменить и добиться такого положения, чтобы $\alpha_1 \neq 0$.

Учитывая разложение вектора $\bar{y}^{(0)}$ по формуле (3.5.2) и принимая во внимание, что $A^k \bar{x}_i = \lambda_i^k \bar{x}_i \quad (i=1, 2, \dots, n)$, получим

$$\bar{y}^{(k)} = A^k \bar{y}^{(0)} = \alpha_1 \lambda_1^k \bar{x}_1 + \alpha_2 \lambda_2^k \bar{x}_2 + \dots + \alpha_n \lambda_n^k \bar{x}_n. \quad (3.5.3)$$

В рассматриваемом степенном методе о собственных значениях и собственных векторах матрицы A судят по последовательности векторов $\bar{y}^{(k)}$. С этой целью введем обозначения:

$$\begin{aligned} \bar{y}^{(k)} &= (y_1^{(k)}, y_2^{(k)}, \dots, y_n^{(k)})' \quad (k=0, 1, 2, \dots), \\ \bar{x}_p &= (x_{1p}, x_{2p}, \dots, x_{np})' \quad (p=1, 2, \dots, n) \end{aligned}$$

и установим связь между компонентами векторов $\bar{y}^{(k)}$ и наибольшим собственным значением λ_1 . Из формулы (3.5.3) получим

$$y_i^{(k)} = \alpha_1 \lambda_1^k x_{i1} + \alpha_2 \lambda_2^k x_{i2} + \dots + \alpha_n \lambda_n^k x_{in} \quad (i=1, 2, \dots, n).$$

Пусть компонента x_{s1} вектора \bar{x}_1 отлична от нуля. В этом случае

$$y_s^{(k)} = \beta_{s1} \lambda_1^k + \beta_{s2} \lambda_2^k + \dots + \beta_{sn} \lambda_n^k, \quad (3.5.4)$$

где $\beta_{si} = \alpha_i x_{si}$ и $\beta_{s1} \neq 0$.

Отношение компонент $y_s^{(k+1)}$ и $y_s^{(k)}$ дает

$$\begin{aligned}\frac{y_s^{(k+1)}}{y_s^{(k)}} &= \frac{\beta_{s1}\lambda_1^{k+1} + \beta_{s2}\lambda_2^{k+1} + \dots + \beta_{sn}\lambda_n^{k+1}}{\beta_{s1}\lambda_1^k + \beta_{s2}\lambda_2^k + \dots + \beta_{sn}\lambda_n^k} = \\ &= \lambda_1 \frac{1 + \gamma_{s2}\mu_2^{k+1} + \dots + \gamma_{sn}\mu_n^{k+1}}{1 + \gamma_{s2}\mu_2^k + \dots + \gamma_{sn}\mu_n^k},\end{aligned}\quad (3.5.5)$$

$$\gamma_{si} = \frac{\beta_{si}}{\beta_{s1}}, \quad \mu_i = \frac{\lambda_i}{\lambda_1}.\quad (3.5.6)$$

Если k достаточно велико, то в силу (3.5.1), (3.5.5) и (3.5.6) получим, что

$$\frac{y_s^{(k+1)}}{y_s^{(k)}} = \lambda_1 + O(|\mu_2|^k),$$

и, следовательно, в качестве λ_1 можно при больших k взять такое приближенное значение:

$$\lambda_1 \approx \frac{y_s^{(k+1)}}{y_s^{(k)}}.\quad (3.5.7)$$

Обычно несколько компонент вектора \bar{x}_1 отличны от нуля, поэтому в формуле (3.5.7) можно вычислять отношения при нескольких значениях s и, если эти отношения в принятой точности вычислений оказываются постоянными, то это означает, что λ_1 вычислено с заданной точностью.

Быстрота сходимости процесса в рассматриваемом случае определяется величиной μ_2 ($|\mu_2| < 1$). Она может быть медленной, если $|\mu_2|$ близок к единице.

При вычислении векторов $\bar{y}^{(k)}$ иногда может оказаться, что компоненты этих векторов быстро растут. Чтобы избежать этого нежелательного явления, можно на каждом шаге нормировать получаемые векторы $\bar{y}^{(k)}$, умножая их, например, на числа $\frac{1}{\|\bar{y}^{(k)}\|_I}$ или $\frac{1}{\|\bar{y}^{(k)}\|_{II}}$. При этом вместо последовательности $\bar{y}^{(k)}$ мы получим последовательность $\bar{z}^{(k)} = \rho_k \bar{y}^{(k)}$, где ρ_k — один из нормирующих множителей. Теперь для получения λ_1 надо брать отношения компонент векторов $A\bar{z}^{(k)}$ и $\bar{z}^{(k)}$.

Рассматриваемый процесс дает возможность определить также собственный вектор матрицы A , отвечающий наибольшему собственному значению λ_1 . Действительно, из (3.5.3) имеем

$$\bar{y}^{(k)} = \lambda_1^k [\alpha_1 \bar{x}_1 + \mu_2^k \alpha_2 \bar{x}_2 + \dots + \mu_n^k \alpha_n \bar{x}_n]. \quad (3.5.8)$$

Если учесть, что $|\mu_i| < 1$ ($i=2, 3, \dots, n$), то из (3.5.8) следует, что при больших k с точностью до постоянного множителя в качестве собственного вектора, отвечающего λ_1 , приближенно можно взять вектор $\bar{y}^{(k)}$. Когда матрица A симметрична, можно легко указать другой вычислительный процесс, более быстро сходящийся к наибольшему по модулю собственному значению λ_1 . Напомним прежде всего, что симметричная матрица A всегда имеет полную систему собственных векторов $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$ и мы вправе эти векторы считать ортонормированными, т. е. такими, что для них выполняется равенство

$$(\bar{x}_i, \bar{x}_j) = \delta_{ij} \quad (i, j = 1, 2, \dots, n).$$

Составим скалярные произведения $(\bar{y}^{(k)}, \bar{y}^{(k)})$ и $(\bar{y}^{(k+1)}, \bar{y}^{(k)})$. С помощью (3.5.3) для них получим

$$(\bar{y}^{(k)}, \bar{y}^{(k)}) = \alpha_1^2 \lambda_1^{2k} + \alpha_2^2 \lambda_2^{2k} + \dots + \alpha_n^2 \lambda_n^{2k},$$

$$(\bar{y}^{(k+1)}, \bar{y}^{(k)}) = \alpha_1^2 \lambda_1^{2k+1} + \alpha_2^2 \lambda_2^{2k+1} + \dots + \alpha_n^2 \lambda_n^{2k+1}$$

и, значит,

$$\frac{(\bar{y}^{(k+1)}, \bar{y}^{(k)})}{(\bar{y}^{(k)}, \bar{y}^{(k)})} = \lambda_1 + O(|\mu_2|^{2k}).$$

Предположим теперь, что собственные значения матрицы A распределены по модулю так:

$$|\lambda_1| > |\lambda_{r+1}| \geq |\lambda_{r+2}| \geq \dots \geq |\lambda_n|$$

и

$$\lambda_1 = \lambda_2 = \dots = \lambda_r.$$

Здесь r означает кратность собственного значения λ_1 . В этом случае формула (3.5.3) верна, но она примет такой вид:

$$\bar{y}^{(k)} = \lambda_1^k (\alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2 + \dots + \alpha_r \bar{x}_r) + \alpha_{r+1} \lambda_{r+1}^k \bar{x}_{r+1} + \dots + \alpha_n \lambda_n^k \bar{x}_n.$$

По аналогии с формулой (3.5.4) можно записать

$$y_s^{(k)} = \beta_{s1}\lambda_1^k + \beta_{s\ r+1}\lambda_{r+1}^k + \dots + \beta_{sn}\lambda_n^k,$$

где

$$\beta_{s1} = \alpha_1 x_{s1} + \alpha_2 x_{s2} + \dots + \alpha_r x_{sr}, \quad \beta_{si} = \alpha_i x_{si} \quad (i = r+1, r+2, \dots, n).$$

Полагаем также, что $\beta_{s1} \neq 0$. Если k достаточно велико, то для вычисления λ_1 получим формулу

$$\frac{y_s^{(k+1)}}{y_s^{(k)}} = \lambda_1 + O(|\mu_{r+1}|^k). \quad (3.5.9)$$

Отметим, что сама формула (3.5.9) не дает возможности судить о кратности собственного значения λ_1 . Как и в предыдущем случае, в качестве собственного вектора матрицы A , соответствующего собственному значению λ_1 , приближенно можно взять вектор $\bar{y}^{(k)}$. Исходя из различных начальных векторов $\bar{y}^{(0)}$, мы, вообще говоря, придем к различным собственным векторам $\bar{y}^{(k)} = A^k \bar{y}^{(0)}$, что даст возможность вычислить другие собственные векторы, отвечающие собственному значению λ_1 .

Рассмотрим также случай, когда матрица A имеет два наибольших по модулю собственных значения и эти значения вещественны и противоположны по знаку. Будем считать, что собственные значения матрицы A распределены по модулю следующим образом:

$$|\lambda_1| = |\lambda_2| > |\lambda_3| \geq |\lambda_4| \geq \dots \geq |\lambda_n|.$$

и

$$\lambda_1 = -\lambda_2.$$

Тогда, в силу формулы (3.5.3), получим

$$\begin{aligned} \bar{y}^{(2k)} &= \alpha_1 \lambda_1^{2k} \bar{x}_1 + \alpha_2 \lambda_2^{2k} \bar{x}_2 + \alpha_3 \lambda_3^{2k} \bar{x}_3 + \dots + \alpha_n \lambda_n^{2k} \bar{x}_n = \\ &= \lambda_1^{2k} (\alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2) + \alpha_3 \lambda_3^{2k} \bar{x}_3 + \dots + \alpha_n \lambda_n^{2k} \bar{x}_n, \\ \bar{y}^{(2k+1)} &= \alpha_1 \lambda_1^{2k+1} \bar{x}_1 + \alpha_2 \lambda_2^{2k+1} \bar{x}_2 + \alpha_3 \lambda_3^{2k+1} \bar{x}_3 + \dots + \alpha_n \lambda_n^{2k+1} \bar{x}_n = \\ &= \lambda_1^{2k+1} (\alpha_1 \bar{x}_1 - \alpha_2 \bar{x}_2) + \alpha_3 \lambda_3^{2k+1} \bar{x}_3 + \dots + \alpha_n \lambda_n^{2k+1} \bar{x}_n. \end{aligned}$$

Отсюда видно, что векторы $\bar{y}^{(2k)}$ и $\bar{y}^{(2k+1)}$ одновременно нельзя использовать для определения λ_1 , ибо у этих векторов различные главные части, а именно: $\lambda_1^{2k} (\alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2)$ — у первого и $\lambda_1^{2k+1} (\alpha_1 \bar{x}_1 - \alpha_2 \bar{x}_2)$ — у второго

вектора. Однако в этом случае мы сможем определить λ_1^2 , используя векторы $\bar{y}^{(2k)}$ и $\bar{y}^{(2k+2)}$ или $\bar{y}^{(2k-1)}$ и $\bar{y}^{(2k+1)}$. Действительно, при некотором p

$$y_s^{(p)} = \beta_{s1}^{(p)} \lambda_1^p + \beta_{s3} \lambda_3^p + \dots + \beta_{sn} \lambda_n^p,$$

где

$$\beta_{s1}^{(p)} = \alpha_1 x_{s1} + (-1)^p \alpha_2 x_{s2}, \quad \beta_{si} = \alpha_i x_{si} \quad (i=3, 4, \dots, n).$$

Если p взять равным $2k$ и $2k+2$, то у компонент $y_s^{(2k)}$ и $y_s^{(2k+2)}$ главные части будут равны соответственно $\lambda_1^{2k} \beta_{s1}^{(2k)}$ и $\lambda_1^2 \lambda_1^{2k} \beta_{s1}^{(2k+2)}$, причем $\beta_{s1}^{(2k)} = \beta_{s1}^{(2k+2)}$. То же самое можно сказать и о компонентах $y_s^{(2k-1)}$ и $y_s^{(2k+1)}$. Следовательно,

$$\frac{y_s^{(2k+2)}}{y_s^{(2k)}} = \lambda_1^2 + O(|\mu_3|^{2k}) \quad \text{и} \quad \frac{y_s^{(2k+1)}}{y_s^{(2k-1)}} = \lambda_1^2 + O(|\mu_3|^{2k-1}).$$

Для нахождения собственных векторов, принадлежащих λ_1 и $\lambda_2 = -\lambda_1$, целесообразно построить векторы

$$\begin{aligned} \bar{y}^{(k)} + \lambda_1 \bar{y}^{(k-1)} &= 2\alpha_1 \lambda_1^{k-1} \bar{x}_1 + \alpha_3 \lambda_3^{k-1} (\lambda_3 + \lambda_1) \bar{x}_3 + \dots + \\ &+ \alpha_n \lambda_n^{k-1} (\lambda_n + \lambda_1) \bar{x}_n = \lambda_1^k [2\alpha_1 \bar{x}_1 + O(|\mu_3|^k)], \\ \bar{y}^{(k)} - \lambda_1 \bar{y}^{(k-1)} &= 2\alpha_2 (-\lambda_1)^{k-1} \bar{x}_2 + \alpha_3 \lambda_3^{k-1} (\lambda_3 - \lambda_1) \bar{x}_3 + \dots + \\ &+ \alpha_n \lambda_n^{k-1} (\lambda_n - \lambda_1) \bar{x}_n = \lambda_2^k [2\alpha_2 \bar{x}_2 + O(|\mu_3|^k)]. \end{aligned}$$

Из этих формул видно, что с точностью до постоянного множителя в качестве собственного вектора, отвечающего λ_1 , можно приближенно взять вектор $\bar{y}^{(k)} + \lambda_1 \bar{y}^{(k-1)}$, а в качестве собственного вектора, отвечающего λ_2 , — вектор $\bar{y}^{(k)} - \lambda_1 \bar{y}^{(k-1)}$.

Если матрица A имеет пару наибольших по модулю комплексно сопряженных собственных значений, то указанные выше приемы нельзя применить к нахождению этих двух собственных значений. Поэтому целесообразно в этом случае несколько видоизменить схему вычислений.

Будем считать, что

$$|\lambda_1| = |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|,$$

$$\lambda_1 = re^{i\theta}, \quad \lambda_2 = re^{-i\theta} \quad \text{и} \quad \lambda_1 = \bar{\lambda}_2.$$

Как и ранее, получим

$$y_s^{(h)} = \beta_{s1}\lambda_1^h + \beta_{s2}\lambda_2^h + \dots + \beta_{sn}\lambda_n^h, \quad (3.5.10)$$

где $\beta_{si} = \alpha_i x_{si}$ ($i=1, 2, \dots, n$). Допустим, что β_{s1} и β_{s2} отличны от нуля. Поскольку матрица A и начальный вектор $y^{(0)}$ вещественны, то вещественным будет и вектор $\bar{y}^{(h)} = A^h \bar{y}^{(0)}$. Значит, в формуле (3.5.10) величины β_{s1} и β_{s2} должны быть комплексно сопряженными. Пусть

$$\beta_{s1} = R_s e^{i\kappa_s}, \quad \beta_{s2} = R_s e^{-i\kappa_s}.$$

Теперь из формулы (3.5.10) получим

$$y_s^{(h)} = 2R_s r^h \cos(k\theta + \kappa_s) + \beta_{s3}\lambda_3^h + \dots + \beta_{sn}\lambda_n^h. \quad (3.5.11)$$

Тот факт, что матрица A имеет пару наибольших по модулю комплексно сопряженных собственных значений, проявляется обычно в сильном колебании по величине и переменам знака в компонентах $y_s^{(h)}$ векторов $\bar{y}^{(k)}$. Так, например, при достаточно больших k и значении аргумента $k\theta + \kappa_s$, близком к нулю, главным в формуле (3.5.11) будет член $2R_s r^h \cos(k\theta + \kappa_s)$, если же аргумент $k\theta + \kappa_s$ близок к $\frac{\pi}{2}$, то главный член выделить, вообще говоря, нельзя. Заметим также, что при изменении k могут быть перемены знака в компонентах $y_s^{(k)}$.

Мы сможем вычислить собственные значения λ_1 и λ_2 , если укажем правило для нахождения чисел r и θ . При достаточно большом k в силу условия $|\lambda_2| > |\lambda_3|$ из формулы (3.5.11) получим

$$y_s^{(k)} = 2R_s r^k \cos(k\theta + \kappa_s) + O(|\lambda_3|^k). \quad (3.5.12)$$

Наряду с равенством (3.5.12) будут иметь место следующие аналогичные равенства:

$$y_s^{(k+1)} = 2R_s r^{k+1} \cos[(k+1)\theta + \kappa_s] + O(|\lambda_3|^{k+1}) \quad (3.5.13)$$

и

$$y_s^{(k+2)} = 2R_s r^{k+2} \cos[(k+2)\theta + \kappa_s] + O(|\lambda_3|^{k+2}). \quad (3.5.14)$$

Эти равенства мы используем для того, чтобы найти приближенные значения r и θ . С этой целью введем в рассмотрение определитель

$$I_s^{(k)} = \begin{vmatrix} y_s^{(k)} & y_s^{(k+1)} \\ y_s^{(k+1)} & y_s^{(k+2)} \end{vmatrix}$$

и вычислим его значение, используя формулы (3.5.12) — (3.5.14). Имеем

$$I_s^{(k)} = 4R_s^2 r^{2k+2} \{ \cos [(k+2)\theta + \kappa_s] \cos (k\theta + \kappa_s) - \cos^2 [(k+1)\theta + \kappa_s] \} + \\ + M_s r^k O(|\lambda_3|^k) = -4R_s^2 r^{2k+2} \sin^2 \theta + M_s r^k O(|\lambda_3|^k),$$

где M_s — некоторая константа. Здесь $\sin \theta \neq 0$, ибо λ_1 и λ_2 — по предположению комплексные числа. Аналогично получим

$$I_s^{(k-1)} = -4R_s^2 r^{2k} \sin^2 \theta + M_s r^{k-1} O(|\lambda_3|^k).$$

Значит, модуль комплексного числа λ_1 можно приближенно вычислить по такой формуле:

$$r^2 \approx \frac{I_s^{(k)}}{I_s^{(k-1)}} \approx \frac{y_s^{(k)} y_s^{(k+2)} - (y_s^{(k+1)})^2}{y_s^{(k-1)} y_s^{(k+1)} - (y_s^{(k)})^2}, \quad (3.5.15)$$

Эти отношения следует определять для нескольких значений s . Совпадение результатов будет свидетельствовать о достижении необходимой точности в вычислении r^2 . После того как мы найдем r , аргумент комплексного числа можно находить приближенно по формуле

$$\cos \theta \approx \frac{y_s^{(k+2)} + r^2 y_s^{(k)}}{2r y_s^{(k+1)}}, \quad (3.5.16)$$

так как

$$y_s^{(k+2)} + r^2 y_s^{(k)} = 2R_s r^{k+2} \{ \cos [(k+2)\theta + \kappa_s] + \cos (k\theta + \kappa_s) \} + \\ + M_s \cdot O(|\lambda_3|^k) = 2R_s r^{k+2} \cos [(k+1)\theta + \kappa_s] \cos \theta + \\ + M_s \cdot O(|\lambda_3|^k) = 2r y_s^{(k+1)} \cos \theta + M_s \cdot O(|\lambda_3|^k).$$

Теперь для λ_1 и λ_2 окончательно получим

$$\lambda_1 = r(\cos \theta + i \sin \theta), \quad \lambda_2 = r(\cos \theta - i \sin \theta).$$

Отметим, что формулы (3.5.15) и (3.5.16) позволяют вычислять величины r^2 и $\cos \theta$ с погрешностью $O\left(\left|\frac{\lambda_3}{\lambda_2}\right|^k\right)$.

Если λ_1 и λ_2 найдены, то легко можно найти и собственные векторы, отвечающие этим собственным значениям. Действительно, в силу формулы (3.5.3) имеем

$$\bar{y}^{(k)} = \alpha_1 \lambda_1^k \bar{x}_1 + \alpha_2 \lambda_2^k \bar{x}_2 + O(|\lambda_3|^k),$$

$$\bar{y}^{(k+1)} = \alpha_1 \lambda_1^{k+1} \bar{x}_1 + \alpha_2 \lambda_2^{k+1} \bar{x}_2 + O(|\lambda_3|^{k+1}).$$

Отсюда находим

$$\bar{y}^{(k+1)} - \lambda_2 \bar{y}^{(k)} = \alpha_1 \lambda_1^k (\lambda_1 - \lambda_2) \bar{x}_1 + O(|\lambda_3|^k),$$

$$\bar{y}^{(k+1)} - \lambda_1 \bar{y}^{(k)} = \alpha_2 \lambda_2^k (\lambda_2 - \lambda_1) \bar{x}_2 + O(|\lambda_3|^k).$$

Таким образом, в качестве собственного вектора, отвечающего собственному значению λ_1 , при больших k приближенно можно взять вектор $\bar{y}^{(k+1)} - \lambda_2 \bar{y}^{(k)}$, другой вектор $\bar{y}^{(k+1)} - \lambda_1 \bar{y}^{(k)}$ в свою очередь можно приближенно принять за собственный вектор, отвечающий собственному значению λ_2 .

Во всех предыдущих случаях вычислительные схемы степенного метода строились в предположении, что матрица A имеет линейные элементарные делители, отвечающие наибольшему по модулю собственному значению. Как это будет показано ниже, в степенном методе имеется возможность вычислять также и наибольшее по модулю собственное значение, которому соответствует нелинейный элементарный делитель. Однако при этом ход итерационного степенного процесса существенно усложняется.

Рассмотрим, например, случай, когда λ_1 вещественно и принадлежит в канонической форме Жордана ящику $\begin{bmatrix} \lambda_1 & 0 \\ 1 & \lambda_1 \end{bmatrix}$. При этом мы будем считать, что другим собственным значениям матрицы A соответствуют линейные элементарные делители и все собственные значения по модулю распределены следующим образом:

$$|\lambda_1| > |\lambda_3| \geq \dots \geq |\lambda_n|.$$

В этом случае при решении задачи о вычислении собственного значения λ_1 удобно использовать вместо базиса из собственных векторов канонический базис. Пусть векторы $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$ образуют канонический базис матрицы A . Известно, что воздействие матрицы A на векторы этого базиса происходит по формулам

$$\left. \begin{aligned} A\bar{x}_1 &= \lambda_1 \bar{x}_1 + \bar{x}_2, \\ A\bar{x}_2 &= \lambda_1 \bar{x}_2, \\ A\bar{x}_3 &= \lambda_3 \bar{x}_3, \\ &\vdots \\ A\bar{x}_n &= \lambda_n \bar{x}_n. \end{aligned} \right\} \quad (3.5.17)$$

Значит,

$$\left. \begin{aligned} A^k \bar{x}_1 &= \lambda_1^k \bar{x}_1 + k \lambda_1^{k-1} \bar{x}_2, \\ A^k \bar{x}_2 &= \lambda_1^k \bar{x}_2, \\ A^k \bar{x}_i &= \lambda_i^k \bar{x}_i \quad (i=3, 4, \dots, n). \end{aligned} \right\} \quad (3.5.18)$$

Покажем теперь, каким образом может быть найдено собственное значение λ_1 при указанном распределении собственных значений матрицы A . Пусть $\bar{y}^{(0)}$ — начальный вектор. Напишем разложение $\bar{y}^{(0)}$ по векторам канонического базиса матрицы A :

$$\bar{y}^{(0)} = \alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2 + \alpha_3 \bar{x}_3 + \dots + \alpha_n \bar{x}_n. \quad (3.5.19)$$

Используя формулы (3.5.18) и (3.5.19), получим

$$\bar{y}^{(k)} = A^{(k)} \bar{y}^{(0)} = \lambda_1^k (\alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2) + \alpha_2 k \lambda_1^{k-1} \bar{x}_2 + \alpha_3 \lambda_3^k \bar{x}_3 + \dots + \alpha_n \lambda_n^k \bar{x}_n. \quad (3.5.20)$$

На основании формулы (3.5.20) можно теперь записать такое выражение для i -й компоненты вектора $\bar{y}^{(k)}$:

$$y_i^{(k)} = \lambda_1^k (\alpha_1 x_{i1} + \alpha_2 x_{i2}) + \alpha_2 k \lambda_1^{k-1} x_{i2} + \alpha_3 \lambda_3^k x_{i3} + \dots + \alpha_n \lambda_n^k x_{in}.$$

Предположим, что при некотором $i=s$ коэффициент

$$\alpha_2 x_{s2} \neq 0.$$

Тогда отношение компонент $y_s^{(k+1)}$ и $y_s^{(k)}$ может быть представлено в виде

$$\frac{y_s^{(k+1)}}{y_s^{(k)}} = \lambda_1 \left(1 + O\left(\frac{1}{k}\right) \right). \quad (3.5.21)$$

Полученная формула показывает, что отношение $\frac{y_s^{(k+1)}}{y_s^{(k)}}$ стремится к λ_1 при $k \rightarrow \infty$.

Однако из-за наличия множителя k во втором слагаемом формулы (3.5.20) сходимость будет медленнее, чем сходимость любой геометрической прогрессии со знаменателем, меньшим единицы. А это означает, что найти в этом случае λ_1 из отношения $\frac{y_s^{(k+1)}}{y_s^{(k)}}$

практически невозможно. Поэтому в рассматриваемом случае целесообразно находить не само собственное значение λ_1 , а коэффициенты $p = -2\lambda_1$ и $q = \lambda_1^2$ квадратного уравнения $\lambda^2 + p\lambda + q = 0$, кратным корнем которого является λ_1 .

Введем обозначения

$$\beta_{s1} = \alpha_1 x_{s1} + \alpha_2 x_{s2}, \quad \beta_{s2} = \alpha_2 x_{s2}, \quad \dots, \quad \beta_{sn} = \alpha_n x_{sn}$$

и запишем выражение для $y_s^{(k)}$:

$$y_s^{(k)} = \beta_{s1} \lambda_1^k + k \beta_{s2} \lambda_1^{k-1} + \beta_{s3} \lambda_3^k + \dots + \beta_{sn} \lambda_n^k.$$

Значит,

$$\begin{aligned} y_s^{(k+1)} + p y_s^{(k)} + q y_s^{(k-1)} &= \beta_{s1} \lambda_1^{k-1} (\lambda_1^2 + p \lambda_1 + q) + \beta_{s2} \lambda_1^{k-2} [k (\lambda_1^2 + p \lambda_1 + q) + (\lambda_1^2 - q)] + \\ &+ O(|\lambda_3|^k) = O(|\lambda_3|^k). \end{aligned}$$

Аналогично при некотором $t \neq s$ ($1 \leq t \leq n$) получим

$$y_t^{(k+1)} + p y_t^{(k)} + q y_t^{(k-1)} = O(|\lambda_3|^k).$$

После того как мы вычислим для рассматриваемой матрицы A векторы $\bar{y}^{(r)}$ и компоненты $y_s^{(r)}, y_t^{(r)}$ ($r = k-1, k, k+1$), искомые значения величин p и q определяются как решение системы

$$y_s^{(k+1)} + p y_s^{(k)} + q y_s^{(k-1)} \approx 0,$$

$$y_t^{(k+1)} + p y_t^{(k)} + q y_t^{(k-1)} \approx 0.$$

Значит,

$$p \approx - \frac{y_s^{(k-1)} y_t^{(k+1)} - y_t^{(k-1)} y_s^{(k+1)}}{y_s^{(k-1)} y_t^{(k)} - y_t^{(k-1)} y_s^{(k)}},$$

$$q \approx \frac{y_s^{(k)} y_t^{(k+1)} - y_t^{(k)} y_s^{(k+1)}}{y_s^{(k-1)} y_t^{(k)} - y_t^{(k-1)} y_s^{(k)}}.$$

Отметим, что эти формулы позволяют определять коэффициенты p и q с точностью до величин порядка $O\left(\left|\frac{\lambda_3}{\lambda_1}\right|^k\right)$.

Для определения собственного значения λ_1 , очевидно, достаточно вычислить один из коэффициентов p или q , ибо $p = -2\lambda_1$, а $q = \lambda_1^2$. В то же время совпадение чисел $-\frac{p}{2}$ и \sqrt{q} , определяемых по найденным значениям p и q , служит контролем правильности предположения о вхождении собственного значения λ_1 в канонический ящик $\begin{bmatrix} \lambda_1 & 0 \\ 1 & \lambda_1 \end{bmatrix}$.

Используя формулу (3.5.20), можно легко найти собственный вектор \bar{x}_2 , отвечающий собственному значению λ_1 . Действительно, в силу названной формулы имеем

$$\left. \begin{aligned} \bar{y}^{(k)} &= \lambda_1^k (\alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2) + \alpha_2 k \lambda_1^{k-1} \bar{x}_2 + O(|\lambda_3|^k), \\ \bar{y}^{(k+1)} &= \lambda_1^{k+1} (\alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2) + \alpha_2 (k+1) \lambda_1^k \bar{x}_2 + O(|\lambda_3|^{k+1}). \end{aligned} \right\} \quad (3.5.22)$$

Отсюда получим

$$\frac{1}{\lambda_1^k} (\bar{y}^{(k+1)} - \lambda_1 \bar{y}^{(k)}) = \alpha_2 \bar{x}_2 + O\left(\left|\frac{\lambda_3}{\lambda_1}\right|^k\right).$$

Таким образом, в качестве собственного вектора, отвечающего собственному значению λ_1 , можно приближенно взять вектор

$$\frac{1}{\lambda_1^k} (\bar{y}^{(k+1)} - \lambda_1 \bar{y}^{(k)}).$$

Формулы (3.5.22) позволяют определить также и вектор $\alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2$. Обозначим $\bar{x} = \alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2$. Тогда

$$\bar{y}^{(k)} = \lambda_1^k \bar{x} + \alpha_2 k \lambda_1^{k-1} \bar{x}_2 + O(|\lambda_3|^k)$$

и

$$\bar{y}^{(k+1)} = \lambda_1^{k+1} \bar{x} + \alpha_2 (k+1) \lambda_1^k \bar{x}_2 + O(|\lambda_3|^{k+1}).$$

Отсюда получим

$$\bar{x} = \frac{1}{\lambda_1^{k+1}} ((k+1) \lambda_1 \bar{y}^{(k)} - k \bar{y}^{(k+1)}) + O\left(k \left| \frac{\lambda_3}{\lambda_1} \right|^k\right).$$

Полученный вектор \bar{x} есть проекция начального вектора $\bar{y}^{(0)}$ на корневое подпространство векторов, соответствующих собственному значению λ_1 . Зная вектор \bar{x} , можно корневой вектор \bar{x}_1 , соответствующий собственному значению λ_1 , определить с точностью до слагаемого, пропорционального собственному вектору \bar{x}_2 , а именно: $\alpha_1 \bar{x}_1 = \bar{x} - \alpha_2 \bar{x}_2$.

Указанный выше порядок вычислений можно применять и тогда, когда собственные значения λ_1 и λ_2 вещественны и выполняются неравенства

$$|\lambda_1| > |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|,$$

причем отношение $\left| \frac{\lambda_2}{\lambda_1} \right|$ близко к единице.

Как это было показано выше, находить λ_1 можно, вообще говоря, по формуле

$$\frac{y_s^{(k+1)}}{y_s^{(k)}} = \lambda_1 + O\left(\left| \frac{\lambda_2}{\lambda_1} \right|^k\right).$$

Однако из-за близости отношения $\left| \frac{\lambda_2}{\lambda_1} \right|$ к единице выражение $\frac{y_s^{(k+1)}}{y_s^{(k)}}$ при возрастании k

будет медленно стремиться к λ_1 , и тем медленнее, чем ближе друг к другу значения $|\lambda_1|$ и $|\lambda_2|$. Поэтому λ_1 и λ_2 следует вычислять по такой схеме. Сначала по заданной матрице A и вектору $\bar{y}^{(0)}$ строим итерационную последовательность векторов $\{\bar{y}^{(k)}\}$. Затем выбираем из этой последовательности векторы $\bar{y}^{(k-1)}$, $\bar{y}^{(k)}$, $\bar{y}^{(k+1)}$ и по компонентам номеров s и t , как и ранее, вычисляем коэффициенты p и q для многочлена $\lambda^2 + p\lambda + q$. Находим корни этого многочлена и принимаем их за искомые собственные значения λ_1 и λ_2 . Заметим, что по смыслу решаемой задачи корни u многочлена должны быть вещественными и близкими по модулю.

3.5.2. Вычисление всех собственных значений положительно определенной симметрической матрицы

Рассмотренные выше случаи позволяют находить для широкого класса вещественных матриц одно или два наибольших по абсолютной величине собственных значения, используя только сведения о распределении значений по абсолютной величине и элементарных делителях матрицы A . При этом в вычислительных схемах не учитываются другие специальные свойства матриц, такие, например, как симметричность, положительная определенность и др.

Оказывается, что в некоторых случаях, используя эти свойства, можно упростить процесс вычислений и получить возможность вычислить все собственные значения и векторы. Так будет, например, в случае, когда вещественная матрица A положительно определенная и симметрическая. Известно, что у этой матрицы все собственные значения вещественны и положительны, собственные векторы $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$ могут быть выбраны так, чтобы удовлетворялось условие ортогональности

$$(\bar{x}_i, \bar{x}_j) = 0 \quad \text{при } i \neq j.$$

Если собственные значения $\lambda_1, \lambda_2, \dots, \lambda_n$ матрицы A занумерованы в порядке невозрастания, то

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0.$$

Напишем систему, из которой определяется собственный вектор \bar{x}_1 , отвечающий собственному значению λ_1 :

[illegible]

Одна из компонент вектора \bar{x}_1 может быть взята произвольной. Положим, например, $x_{n1}=1$. Тогда система (3.5.23) будет системой n нелинейных алгебраических уравнений с n неизвестными $x_{11}, x_{21}, \dots, x_{n-1,1}, \lambda_1$. Приведем эту систему к виду, удобному для применения метода итерации:

$$\left. \begin{aligned} x_{11} &= \frac{1}{\lambda_1} (a_{11}x_{11} + a_{12}x_{21} + \dots + a_{1n}), \\ x_{21} &= \frac{1}{\lambda_1} (a_{21}x_{11} + a_{22}x_{21} + \dots + a_{2n}), \\ &\vdots \\ x_{n-1,1} &= \frac{1}{\lambda_1} (a_{n-1,1}x_{11} + a_{n-1,2}x_{21} + \dots + a_{n-1,n}), \\ \lambda_1 &= a_{n1}x_{11} + a_{n2}x_{21} + \dots + a_{nn}. \end{aligned} \right\} \quad (3.5.24)$$

Систему (3.5.24) можно пытаться решать методом итераций, выбирая некоторые начальные приближения для компонент собственного вектора и собственного значения $x_{11}^{(0)}, x_{21}^{(0)}, \dots, x_{n-1,1}^{(0)}$ и $\lambda_1^{(0)}$. Вычислительные формулы будут иметь следующий вид:

$$x_{i1}^{(k+1)} = \frac{1}{\lambda_1^{(k)}} \left(\sum_{j=1}^{n-1} a_{ij} x_{j1}^{(k)} + a_{in} \right) \quad (i=1, 2, \dots, n-1),$$

$$\lambda_1^{(k+1)} = \sum_{j=1}^{n-1} a_{nj} x_{j1}^{(k)} + a_{nn} \quad (k=0, 1, 2, \dots).$$

Заметим, что при вычислении $\lambda_1^{(k+1)}$ вместо значений $x_{j1}^{(k)}$ можно брать значения $x_{j1}^{(k+1)}$. Можно также для решения системы (3.5.24) использовать метод Зейделя.

Если метод итерации для системы (3.5.24) при указанном начальном приближении сходится, то для достаточно больших значений k можно принять приближенно

$$\lambda_1 \approx \lambda_1^{(k)}$$

и

$$\bar{x}_1 \approx (x_{11}^{(k)}, x_{21}^{(k)}, \dots, x_{n-1,1}^{(k)}, 1)'$$

Чтобы найти второе собственное значение λ_2 и соответствующий ему собственный вектор \bar{x}_2 , воспользуемся опять системой, порождающей λ_2 и \bar{x}_2 . Запишем эту систему в виде

$$\lambda_2 x_{i2} = \sum_{j=1}^n a_{ij} x_{j2} \quad (i=1, 2, \dots, n). \quad (3.5.25)$$

Используя условие ортогональности векторов \bar{x}_1 и \bar{x}_2 , получим

$$(\bar{x}_1, \bar{x}_2) = x_{11}^{(k)} x_{12} + x_{21}^{(k)} x_{22} + \dots + x_{n-1,1}^{(k)} x_{n-1,2} + x_{n2} = 0.$$

Выразим отсюда, например, x_{n2} через другие компоненты вектора \bar{x}_2 и найденное выражение для x_{n2} подставим в систему (3.5.25). При этом условии система (3.5.25) может быть записана в эквивалентной форме

$$x_{i2} = \frac{1}{\lambda_2} \sum_{j=1}^{n-1} a_{ij}^{(1)} x_{j2} \quad (i=1, 2, \dots, n-2),$$

$$\lambda_2 = \frac{1}{x_{n-1,2}} \sum_{j=1}^{n-1} a_{n-1,j}^{(1)} x_{j2},$$

где $a_{ij}^{(1)} = a_{ij} - x_{ji}^{(k)} a_{in}$.

Положив $x_{n-1\ 2}=1$ и выбрав начальное приближение $x_{12}^{(0)}, x_{22}^{(0)}, \dots, x_{n-2\ 2}^{(0)}, \lambda_2^{(0)}$, решаем эту систему опять по методу итерации. При условии сходимости итерационной последовательности в качестве λ_2 и x_2 можно принять приближения некоторого номера k

$$\lambda_2 \approx \lambda_2^{(k)}$$

и

$$\bar{x}_2 \approx (x_{12}^{(k)}, x_{22}^{(k)}, \dots, x_{n-2\ 2}^{(k)}, 1, x_{n2})'.$$

Последнее уравнение системы (3.5.25) при $i=n$ можно использовать для контроля правильности вычисления $\lambda_2^{(k)}$ и компонент $x_{12}^{(k)}, x_{22}^{(k)}, \dots, x_{n-2\ 2}^{(k)}$. Это уравнение должно выполняться с необходимой точностью при подстановке в него указанных значений $\lambda_2^{(k)}$ и $x_{i2}^{(k)}$, а также значений $x_{n-1\ 2}=1$ и $x_{n2} = -\sum_{j=1}^{n-2} x_{j1}^{(k)} x_{j2}^{(k)} - x_{n-1\ 1}^{(k)}$.

Аналогично определяются другие собственные значения λ_j ($j=3, 4, \dots, n$) и соответствующие им собственные векторы \bar{x}_j . Следует отметить, что, в силу изложенной здесь схемы метода, последующие собственные значения и векторы могут быть вычислены, вообще говоря, с меньшей точностью, чем предыдущие.

Рассматриваемый метод может иметь исключительные случаи, связанные с тем, что при определении собственного значения λ_s компонента $x_{n-s+1\ s}$ собственного вектора \bar{x}_s может оказаться равной нулю, например уже на первом шаге будет особый случай, если $x_{n1}=0$. Поскольку такая особенность варианта возникает из-за избранного способа приведения системы (3.5.23) к виду, удобному для применения метода итерации, то ее можно устранить. Но мы на этом вопросе останавливаться не будем.

3.5.3. Видоизменения степенного метода

Изложенный в п. 3.5.1 степенной метод может быть усовершенствован в смысле ускорения сходимости получающихся итерационных последовательностей. В основе одного такого видоизменения, пригодного для симметрических матриц, лежит идея сдвига собственных значений таким образом, чтобы величина одного из них стала достаточно малой. При этом величина, обратная сдвинутому собственному значению, будет большой и составит при итерациях вектора $\bar{y}^{(0)}$ обратной матрицей A^{-1} главную часть в обычных для степенного метода разложениях. Это может

быть использовано затем для вычисления непосредственно соответствующего собственного значения. Пусть A — симметрическая матрица и $\lambda_1, \lambda_2, \dots, \lambda_n$ — ее собственные значения. Будем называть числа Λ_k сдвинутыми на величину μ собственными значениями матрицы A , если эти числа являются собственными значениями матрицы $A - \mu E$, т. е. если они удовлетворяют равенствам

$$\Lambda_k = \lambda_k - \mu \quad (k = 1, 2, \dots, n).$$

Рассмотрим некоторый начальный вектор $\bar{y}^{(0)}$ и построим итерационные последовательности $\{\mu_k\}$ и $\{\bar{y}^{(k)}\}$, каждый член которых определяется по формулам

$$\mu_k = \frac{(A\bar{y}^{(k-1)}, \bar{y}^{(k-1)})}{(\bar{y}^{(k-1)}, \bar{y}^{(k-1)})}, \quad (A - \mu_k E)\bar{y}^{(k)} = \bar{y}^{(0)} \quad (3.5.26)$$

$$(k = 1, 2, \dots).$$

Выясним смысл величин μ_k и свойства последовательности $\{\mu_k\}$. С этой целью запишем разложение вектора $\bar{y}^{(0)}$ по собственным векторам матрицы A , отвечающим попарно различным собственным значениям $\lambda_1, \lambda_2, \dots, \lambda_s$, а именно положим, что

$$\bar{y}^{(0)} = \alpha_1 \bar{x}_1 + \alpha_2 \bar{x}_2 + \dots + \alpha_s \bar{x}_s, \quad (3.5.27)$$

где векторы \bar{x}_i удовлетворяют условиям

$$(\bar{x}_i, \bar{x}_j) = \delta_{ij} \quad (i, j = 1, 2, \dots, s).$$

Используя формулы (3.5.26) и (3.5.27), получим

$$\bar{y}^{(k)} = (A - \mu_k E)^{-1} \bar{y}^{(0)} = \frac{\alpha_1}{\lambda_1 - \mu_k} \bar{x}_1 + \frac{\alpha_2}{\lambda_2 - \mu_k} \bar{x}_2 + \dots + \frac{\alpha_s}{\lambda_s - \mu_k} \bar{x}_s,$$

ибо матрица $(A - \mu_k E)^{-1}$ воздействует на векторы \bar{x}_i по формулам

$$(A - \mu_k E)^{-1} \bar{x}_i = \frac{1}{\lambda_i - \mu_k} \bar{x}_i \quad (i = 1, 2, \dots, s).$$

Следовательно, для μ_{k+1} получим

$$\begin{aligned}\mu_{k+1} &= \frac{(A\bar{y}^{(k)}, \bar{y}^{(k)})}{(\bar{y}^{(k)}, \bar{y}^{(k)})} = \\ &= \frac{\frac{\alpha_1^2}{(\lambda_1 - \mu_k)^2} \lambda_1 + \frac{\alpha_2^2}{(\lambda_2 - \mu_k)^2} \lambda_2 + \dots + \frac{\alpha_s^2}{(\lambda_s - \mu_k)^2} \lambda_s}{\frac{\alpha_1^2}{(\lambda_1 - \mu_k)^2} + \frac{\alpha_2^2}{(\lambda_2 - \mu_k)^2} + \dots + \frac{\alpha_s^2}{(\lambda_s - \mu_k)^2}}.\end{aligned}$$

Если, например, сдвинутое на величину μ_k собственное значение $\Lambda_1 = \lambda_1 - \mu_k$ мало, то в предыдущей формуле в числителе и знаменателе выделится главный член за счет дроби $\frac{\alpha_1^2}{(\lambda_1 - \mu_k)^2}$. Значит, главная часть величины μ_{k+1} должна быть близкой к λ_1 . Рассмотрим разность

$$\mu_{k+1} - \lambda_1 = \frac{\frac{\alpha_2^2}{(\lambda_2 - \mu_k)^2} (\lambda_2 - \lambda_1) + \dots + \frac{\alpha_s^2}{(\lambda_s - \mu_k)^2} (\lambda_s - \lambda_1)}{\frac{\alpha_1^2}{(\lambda_1 - \mu_k)^2} + \dots + \frac{\alpha_s^2}{(\lambda_s - \mu_k)^2}}.$$

Отсюда следует, что

$$\begin{aligned}|\mu_{k+1} - \lambda_1| &= (\mu_k - \lambda_1)^2 \frac{\frac{\alpha_2^2}{(\lambda_2 - \mu_k)^2} |\lambda_2 - \lambda_1| + \dots + \frac{\alpha_s^2}{(\lambda_s - \mu_k)^2} |\lambda_s - \lambda_1|}{\alpha_1^2 + (\mu_k - \lambda_1)^2 \left[\frac{\alpha_2^2}{(\lambda_2 - \mu_k)^2} + \dots + \frac{\alpha_s^2}{(\lambda_s - \mu_k)^2} \right]} = \\ &= O \left((\mu_k - \lambda_1)^2 \right).\end{aligned}\quad (3.5.28)$$

На основании формулы (3.5.28) можно записать, что

$$|\mu_{k+1} - \lambda_1| \leq q (\mu_k - \lambda_1)^2,$$

где $q > 0$ — некоторая константа. Отсюда получим

$$\begin{aligned}|\mu_{k+1} - \lambda_1| &\leq q (\mu_k - \lambda_1)^2 \leq q [q (\mu_{k-1} - \lambda_1)^2]^2 \leq \dots \leq \\ &\leq q^{1+2+\dots+2^{k-1}} (\mu_1 - \lambda_1)^{2^k} = q^{2^k - 1} (\mu_1 - \lambda_1)^{2^k} = \frac{1}{q} [q (\mu_1 - \lambda_1)]^{2^k}.\end{aligned}\quad (3.5.29)$$

Оценка (3.5.29) показывает, что если $\bar{y}^{(0)}$ выбрано удачно, т. е. $q (\mu_1 - \lambda_1) < 1$, то последовательность $\{\mu_k\}$ сходящаяся и

$$\lim_{k \rightarrow \infty} \mu_k = \lambda_1.$$

Из формулы (3.5.29) видно, что сходимость квадратическая. В аналогичном случае степенной метод сходится со скоростью геометрической прогрессии, т. е. гораздо медленнее, чем предлагаемое видоизменение метода.

Из формулы

$$\bar{y}^{(k)} = \frac{\alpha_1}{\lambda_1 - \mu_k} \bar{x}_1 + \frac{\alpha_2}{\lambda_2 - \mu_k} \bar{x}_2 + \dots + \frac{\alpha_s}{\lambda_s - \mu_k} \bar{x}_s$$

видно, что при достаточно больших k и при условии сходимости последовательности $\{\mu_k\}$ в этом случае в качестве собственного вектора, отвечающего собственному значению λ_1 , с точностью до произвольного постоянного множителя приближенно можно взять вектор $\bar{y}^{(k)}$.

Практически метод реализуется так. Сначала вычисляем число $\mu_1 = \frac{(A\bar{y}^{(0)}, \bar{y}^{(0)})}{(\bar{y}^{(0)}, \bar{y}^{(0)})}$, затем, решая систему $(A - \mu_1 E)\bar{y}^{(1)} = \bar{y}^{(0)}$, находим вектор $\bar{y}^{(1)} = (y_1^{(1)}, y_2^{(1)}, \dots, y_n^{(1)})'$, потом, аналогично, по формулам (3.5.26) вычисляем μ_2 и $\bar{y}^{(2)}$ и продолжаем процесс до тех пор, пока два последовательных приближения μ_k и μ_{k+1} совпадут друг с другом на заданное число знаков.

В рассматриваемом выше итерационном процессе векторы $\bar{y}^{(k)}$ определялись из формулы $(A - \mu_k E)\bar{y}^{(k)} = \bar{y}^{(0)}$, в которой правая часть фиксирована и равна $\bar{y}^{(0)}$. Если правую часть менять на каждом шаге и вычислять векторы $\bar{y}^{(k)}$ по формуле $(A - \mu_k E)\bar{y}^{(k)} = \bar{y}^{(k-1)}$ ($k=1, 2, \dots$), а числа μ_k по формуле (3.5.26), то при симметрической матрице A получающийся итерационный процесс будет иметь кубическую сходимость [9].

Укажем еще на один прием, позволяющий иногда быстрее, чем это имеет место в степенном методе, получать искомый результат при решении задачи о вычислении наибольшего по абсолютной величине собственного значения матрицы A . Этот прием связан с последовательным получением четных высоких степеней матрицы A , а именно: $A^2, A^4, A^8, A^{16}, \dots$. На первый взгляд, такое вычисление степеней матриц может быть сопряжено с увеличением объема работы по сравнению со степенным методом, ибо, например, возведение матрицы A в квадрат по объему работы равносильно образованию n итераций вектора $\bar{y}^{(0)}$ матрицей A . Вычисление же матрицы A^{2^k} по указанному выше правилу, следовательно, равносильно в смысле числа выполняемых операций построению kn итераций вектора $\bar{y}^{(0)}$ матрицей A . Таким образом, при равном объеме работы мы, используя указанный прием, сможем, например, вычислить $\bar{y}^{(2^k)} = A^{2^k} \bar{y}^{(0)}$, а в степенном методе сможем вычислить $\bar{y}^{(k^{n+1})}$. Если при этом окажется, что

$kn+1 < 2^k$, то это будет означать, что при одинаковом числе операций мы сможем по видоизмененному методу вычислить более старшие члены итерационной последовательности по сравнению с последним итерационным членом $\bar{y}^{(kn+1)}$, вычисляемым в этом случае по степенному методу. В этом и заключается преимущество перед степенным методом указанного выше видоизменения.

Отметим, что при нахождении итераций вектора порядок вычислений может быть различным. Можно, например, ограничиться вычислением некоторой фиксированной степени матрицы A , а затем составлять итерации посредством вычисленной степени матрицы. Так, если нам нужно вычислить $\bar{y}^{(49)} = A^{49}\bar{y}^{(0)}$, то мы, вычислив, например, A^2, A^4, A^8, A^{16} , найдем затем $A^{16}\bar{y}^{(0)}, A^{16}(A^{16}\bar{y}^{(0)}), A^{16}[A^{16}(A^{16}\bar{y}^{(0)})]$ и, наконец,

$$A\{A^{16}[A^{16}(A^{16}\bar{y}^{(0)})]\}.$$

Если мы вычислили вектор $\bar{y}^{(2^k)} = A^{2^k}\bar{y}^{(0)}$, то далее наибольшее по абсолютной величине собственное значение матрицы A следует вычислять по одному из правил, изложенных в п. 3.5.1, в зависимости от того, как расположены по абсолютной величине собственные значения матрицы и какова ее жорданова форма.

Если собственные значения матрицы A по абсолютной величине расположены следующим образом:

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|, \quad (3.5.30)$$

то процесс вычисления степеней матрицы можно использовать непосредственно для определения наибольшего по модулю собственного значения λ_1 , которое мы будем считать вещественным.

Образую последовательность матриц

$$A, A^2, A^4, A^8, \dots, A^{2^k}, \dots$$

Известно, что

$$\sum_{i=1}^n \lambda_i = \text{Sp } A, \quad \sum_{i=1}^n \lambda_i^2 = \text{Sp } A^2, \dots, \sum_{i=1}^n \lambda_i^m = \text{Sp } A^m.$$

Будем считать, что $m = 2^k$. Значит,

$$\lambda_1^m + \lambda_2^m + \dots + \lambda_n^m = \lambda_1^m \left[1 + \left(\frac{\lambda_2}{\lambda_1} \right)^m + \dots + \left(\frac{\lambda_n}{\lambda_1} \right)^m \right] = \text{Sp } A^m.$$

Отсюда следует, что

$$|\lambda_1| \left[1 + \left(\frac{\lambda_2}{\lambda_1} \right)^m + \dots + \left(\frac{\lambda_n}{\lambda_1} \right)^m \right]^{\frac{1}{m}} = \sqrt[m]{\text{Sp } A^m}. \quad (3.5.31)$$

Учитывая неравенства (3.5.30), из формулы (3.5.31) получим

$$|\lambda_1| = \lim_{m \rightarrow \infty} \sqrt[m]{\text{Sp } A^m}.$$

Таким образом, при достаточно большом m мы можем положить

$$|\lambda_1| \approx \sqrt[m]{\text{Sp } A^m}. \quad (3.5.32)$$

Формула (3.5.32) неудобна тем, что при больших n она требует извлечения корня высокой степени. Поэтому выгоднее пользоваться следующей очевидной формулой:

$$\lambda_1 \approx \frac{\text{Sp } A^{2^k+1}}{\text{Sp } A^{2^k}}.$$

При этом нет необходимости вычислять всю матрицу A^{2^k+1} . Достаточно определить ее след, т. е. определить диагональные элементы и их сумму.

На этом мы заканчиваем изложение теории степенного метода и некоторых его видоизменений, предназначенных в основном для вычисления наибольшего по модулю собственного значения матрицы.

3.5.4. Метод λ -разности

Этот метод позволяет находить собственное значение λ_2 после вычисления λ_1 при условии, что

$$|\lambda_1| > |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|. \quad (3.5.33)$$

Рассмотрим некоторый начальный вектор $\bar{y}^{(0)}$ и предположим, что для него имеет место разложение (3.5.2) по собственным векторам матрицы A , причем будем считать, что в этом разложении коэффициенты α_1 и α_2 отличны от нуля.

Пусть мы вычислили последовательность

$$\bar{y}^{(0)}, \bar{y}^{(1)}, \dots, \bar{y}^{(m)}, \dots, \bar{y}^{(k)}, \dots,$$

где $A^k \bar{y}^{(0)} = \bar{y}^{(k)}$, и по какому-либо методу вычислили λ_1 .

Введем обозначение

$$\Delta_\lambda \bar{y}^{(k)} = \bar{y}^{(k+1)} - \lambda \bar{y}^{(k)} \quad (k=0, 1, 2, \dots). \quad (3.5.34)$$

Величину $\Delta_\lambda \bar{y}^{(k)}$ будем называть λ -разностью от $\bar{y}^{(k)}$. На основании формул (3.5.3) и (3.5.34) для некоторой компоненты s получим

$$\Delta_{\lambda_1} y_s^{(k)} = y_s^{(k+1)} - \lambda_1 y_s^{(k)} = \alpha_2 (\lambda_2 - \lambda_1) \lambda_2^k x_{s2} + \dots + \alpha_n (\lambda_n - \lambda_1) \lambda_n^k x_{sn}$$

и

$$\Delta_{\lambda_1} y_s^{(k-1)} = y_s^{(k)} - \lambda_1 y_s^{(k-1)} = \alpha_2 (\lambda_2 - \lambda_1) \lambda_2^{k-1} x_{s2} + \dots + \alpha_n (\lambda_n - \lambda_1) \lambda_n^{k-1} x_{sn}.$$

Если k достаточно велико, то в выражениях для $\Delta_{\lambda_1} y_s^{(k)}$ и $\Delta_{\lambda_1} y_s^{(k-1)}$ преобладающими будут члены, содержащие λ_2^k . Значит,

$$\lambda_2 \approx \frac{\Delta_{\lambda_1} y_s^{(k)}}{\Delta_{\lambda_1} y_s^{(k-1)}}. \quad (3.5.35)$$

Заметим, что если мы определяли λ_1 по формуле

$$\lambda_1 \approx \frac{y_s^{(k+1)}}{y_s^{(k)}},$$

то λ_2 целесообразно определять по следующей формуле:

$$\lambda_2 \approx \frac{\Delta_{\lambda_1} y_s^{(m)}}{\Delta_{\lambda_1} y_s^{(m-1)}}, \quad (3.5.36)$$

где $m < k$ и m является наименьшим из чисел, при котором преобладание λ_2^m над следующими членами λ_s^m ($s=3, 4, \dots, n$) уже начинает сказываться. Эта формула имеет преимущество перед формулой (3.5.35) в том, что здесь при определении $\Delta_{\lambda_1} y_s^{(m)}$ и $\Delta_{\lambda_1} y_s^{(m-1)}$ нам не приходится вычитать близкие друг к другу числа (имеется в виду $y_s^{(m+1)}$, $\lambda_1 y_s^{(m)}$ и $y_s^{(m)}$, $\lambda_1 y_s^{(m-1)}$), в то время как в случае формулы (3.5.35) мы имеем

$$y_s^{(k+1)} - \lambda_1 y_s^{(k)} \approx 0 \quad \text{и} \quad y_s^{(k)} - \lambda_1 y_s^{(k-1)} \approx 0.$$

В качестве собственного вектора матрицы A , отвечающего λ_2 , приближенно можно взять вектор $\Delta_{\lambda_1} \bar{y}^{(k)}$, ибо

$$\Delta_{\lambda_1} \bar{y}^{(k)} = \alpha_2 (\lambda_2 - \lambda_1) \lambda_2^k \bar{x}_2 + \dots + \alpha_n (\lambda_n - \lambda_1) \lambda_n^k \bar{x}_n.$$

Теоретически возможно метод λ -разности применять и к вычислению следующих собственных значений, однако результаты будут еще менее надежными, чем в случае λ_2 . Причина этого явления кроется в том, что названные вычисления связаны с операцией уничтожения главной части в линейных выражениях вида (3.5.3). А это влечет за собой большую потерю значащих цифр.

В заключение следует сказать, что многим изложенным здесь итерационным методам нахождения собственных значений и собственных векторов матриц присущи наряду с такими положительными свойствами, как простота вычислительного алгоритма, возможность контроля точности получаемого результата и др., и существенные недостатки. Это в первую очередь их медленная сходимость при определении наибольшего по модулю собственного значения и еще более медленная сходимость при определении последующих значений.

В связи с высказанным выше замечанием приобретают важное значение методы, предназначенные для ускорения сходимости таких итерационных процессов. Теория некоторых названных методов будет изложена в § 3.7.

§ 3.6. МЕТОД ВРАЩЕНИЙ

Элементарные унитарные матрицы и матрицы вращения, рассмотренные нами в гл. 2 при решении систем линейных алгебраических уравнений, можно использовать также для построения итерационных процессов, решающих полную проблему собственных значений для симметрических и эрмитовых матриц. Один из таких процессов, предложенный Якоби, известен еще с середины прошлого века. Однако он долгое время не находил практического применения из-за большого объема вычислений, необходимых для его реализации. И лишь с появлением быстродействующих электронных вычислительных машин стало возможным его широкое применение, которое показало, что метод вращений является одним из самых эффективных методов решения полной проблемы собственных значений симметрических и эрмитовых матриц.

Между прочим, пример с историей метода вращений еще раз показывает необходимость и большую значимость в вычислительной математике, равно как и в других областях науки, теоретических исследований, практическая реализация которых на первых порах затруднительна или вообще невозможна.

Отметим, что в настоящее время известно большое число итерационных процессов, предназначенных для решения полной проблемы собственных значений симметрических и эрмитовых матриц.

Основой для построения таких процессов служит известная теорема из алгебры, утверждающая, что если A — эрмитова матрица, то существует такая унитарная матрица V , что преобразование подобия с этой матрицей приводит A к диагональному виду, т. е.

$$V^{-1}AV = \Lambda, \quad (3.6.1)$$

где Λ — диагональная матрица из собственных значений матрицы A . Так как для унитарной матрицы выполняется условие $V^*V = E$, то $V^{-1} = V^*$,

и, значит, формулу (3.6.1) можно записать иначе:

$$V^*AV = \Lambda. \quad (3.6.2)$$

Равенство (3.6.2) не может быть использовано для прямого вычисления элементов матрицы V и диагональной матрицы Λ , ибо оно представляет собою, вообще говоря, систему n^2 уравнений с $n^2 + n$ неизвестными (n^2 элементов матрицы V плюс n элементов матрицы Λ). Однако имеется возможность трактовать задачу приведения заданной эрмитовой матрицы A к диагональному виду как приближенную задачу в следующем смысле.

Предположим, что мы преобразованиями типа (3.6.2) привели матрицу A к некоторой матрице $\tilde{\Lambda}$ вида

$$\tilde{\Lambda} = \begin{bmatrix} \tilde{\lambda}_1 & \lambda_{12} & \dots & \lambda_{1n} \\ \lambda_{21} & \tilde{\lambda}_2 & \dots & \lambda_{2n} \\ \dots & \dots & \dots & \dots \\ \lambda_{n1} & \lambda_{n2} & \dots & \tilde{\lambda}_n \end{bmatrix}. \quad (3.6.3)$$

Предположим также, что внедиагональные элементы матрицы таковы, что величинами

$$\sigma_1 = \sum_{i=2}^n |\lambda_{1i}|^2, \sigma_2 = \sum_{i=1, i \neq 2}^n |\lambda_{2i}|^2, \dots, \sigma_{n-1} = \sum_{i=1, i \neq n-1}^n |\lambda_{n-1, i}|^2, \sigma_n = \sum_{i=1}^{n-1} |\lambda_{ni}|^2$$

по сравнению соответственно с $\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_n$ можно пренебречь. Тогда мы получим

$$\tilde{V}^*A\tilde{V} = \tilde{\Lambda}$$

и, значит, в силу разложения (3.6.2), в качестве приближенных собственных значений матрицы A можно в пределах принятой точности взять числа $\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_n$. Эти числа, естественно, совпадут с точными собственными значениями $\lambda_1, \lambda_2, \dots, \lambda_n$ матрицы A , если все $\sigma_i = 0$ ($i = 1, 2, \dots, n$).

Высказанные выше соображения можно использовать при построении итерационного процесса для вычисления всех собственных значений эрмитовой матрицы A . Действительно, пусть $V_{ij}(\varphi, \psi)$ — некоторые элементарные унитарные матрицы. О правиле построения таких матриц мы еще будем говорить (см. формулу (3.6.32)). Подвергнем с помощью упомянутых матриц преобразования подобия матрицу A и предположим, что в результате мы получили последовательность матриц

$$A^{(0)} = A, A^{(1)}, A^{(2)}, \dots, A^{(k)}, \dots \quad (3.6.4)$$

Будем считать, что каждый элементарный шаг преобразований, заключающийся в умножении матрицы A слева на $V_{ij}^*(\varphi, \psi)$ и справа на $V_{ij}(\varphi, \psi)$, выбирается таким образом, что матрица $A^{(k)}$ при $k \rightarrow \infty$ сколь угодно близко приближается к диагональной матрице. Ниже мы покажем, что такой выбор матриц $V_{ij}(\varphi, \psi)$ возможен.

Близость эрмитовой матрицы A к диагональной мы будем характеризовать числом

$$t^2(A) = \sum_{i, j=1, i \neq j}^n |a_{ij}|^2,$$

т. е. суммой квадратов модулей всех недиагональных элементов матрицы A . Эта близость может быть охарактеризована также любой нормой матрицы $A - D$, где D — диагональная матрица, составленная из диагональных элементов матрицы A .

Итерационный процесс построения последовательности (3.6.4) будем называть монотонным, если выполняется условие

$$t^2(A^{(k)}) < t^2(A^{(k-1)}) \quad (k=1, 2, \dots).$$

Отметим еще, что если разложение вида (3.6.2) найдено, то легко могут быть указаны правила для вычисления собственных векторов матрицы A . Пусть λ_i — i -й диагональный элемент матрицы Λ и $\bar{e}_i = (0, 0, \dots, 0, 1, 0, \dots, 0)^{(i)}$ — соответствующий этому собственному значению собственный вектор матрицы Λ . Тогда

$$\Lambda \bar{e}_i = V^* A V \bar{e}_i = \lambda_i \bar{e}_i$$

или

$$A V \bar{e}_i = \lambda_i V \bar{e}_i. \quad (3.6.5)$$

Обозначим через

$$\bar{x}_i = V \bar{e}_i = (u_{1i}, u_{2i}, \dots, u_{ni})'$$

i -й столбец матрицы V . Теперь формула (3.6.5) примет вид

$$A \bar{x}_i = \lambda_i \bar{x}_i$$

и, значит, в качестве собственного вектора матрицы A , отвечающего собственному значению λ_i , можно взять вектор \bar{x}_i , компонентами которого являются элементы i -го столбца матрицы V .

Приступим теперь к построению формул итерационного процесса. В целях упрощения выкладок будем сначала рассматривать вещественные симметрические матрицы и для них подробно изложим теорию метода вращений, а затем укажем, как полученные результаты переносятся на эрмитовы матрицы.

3.6.1. Случай вещественных симметрических матриц

Пусть A — одна из названных матриц. Для такой матрицы метод вращений заключается в построении последовательности матриц

$$A^{(0)} = A, A^{(1)}, A^{(2)}, \dots, A^{(k)}, \dots,$$

в которой каждая последующая матрица получается из предыдущей при помощи элементарного шага, состоящего в преобразовании подобия предыдущей матрицы посредством некоторой ортогональной матрицы вращения $V_{ij}(\varphi)$ вида

$$V_{ij}(\varphi) = \begin{bmatrix} 1 & & & & & & 0 \\ & \ddots & & & & & \\ & & 1 & & & & \\ & & & \cos \varphi & \dots & -\sin \varphi & \\ & & & \vdots & 1 & \vdots & \\ & & & \vdots & \ddots & \vdots & \\ & & & \sin \varphi & \dots & \cos \varphi & \\ 0 & & & & & & 1 & \ddots & \\ & & & & & & & \ddots & 1 \end{bmatrix} \quad \begin{matrix} (i) \\ \\ \\ (j) \end{matrix} \quad (3.6.6)$$

Пусть указанные выше преобразования доведены до k -го шага, и мы получили матрицу $A^{(k)} = (a_{rs}^{(k)})$. Построим формулы для определения из матрицы $A^{(k)}$ следующей матрицы $A^{(k+1)}$. Поскольку конечной целью итерационного процесса является диагонализация матрицы A , то матрицы $V_{ij}(\varphi)$ мы должны будем выбирать таким образом, чтобы образующийся процесс был монотонным, т. е. чтобы выполнялось условие

$$t^2(A^{(k+1)}) < t^2(A^{(k)}).$$

Существует много способов выбора матрицы $V_{ij}(\varphi)$, при которых это условие будет выполняться. Мы остановимся на способе, который быстрее всех других способов приводит к цели. А именно, рассмотрим матрицу $A^{(k)}$ и найдем в ней максимальный по абсолютной величине внедиагональный элемент. Поскольку $A^{(k)}$ — симметрическая матрица, то условимся считать искомым некоторый элемент $a_{ij}^{(k)}$, для которого $i < j$. На основании этих индексов i, j и элемента $a_{ij}^{(k)}$ построим ортогональную матрицу $V_{ij}^{(k)}$ по формуле

$$V_{ij}^{(k)} = V_{ij}(\varphi^{(k)}),$$

где матрица $V_{ij}(\varphi^{(k)})$ получается из матрицы $V_{ij}(\varphi)$ заменой φ на $\varphi^{(k)}$. Значение параметра $\varphi^{(k)}$ нам еще предстоит определить. Будем находить $\varphi^{(k)}$ из условия равенства нулю у матрицы

$$A^{(k+1)} = V_{ij}^{(k)'} A^{(k)} V_{ij}^{(k)} \quad (3.6.7)$$

элемента $a_{ij}^{(k+1)}$. Укажем формулы для вычисления элементов матрицы $A^{(k+1)}$. Обозначим

$$B^{(k)} = A^{(k)} V_{ij}^{(k)} \quad \text{и} \quad B^{(k)} = (b_{rs}^{(k)}).$$

Тогда, в силу формулы (3.6.6), матрица $B^{(k)}$ во всех столбцах, кроме i -го и j -го, будет иметь те же элементы, что и матрица $A^{(k)}$. Для элементов i -го и j -го столбцов имеют место соответственно формулы

$$\left. \begin{aligned} b_{vi}^{(k)} &= a_{vi}^{(k)} \cos \varphi^{(k)} + a_{vj}^{(k)} \sin \varphi^{(k)}, \\ b_{vj}^{(k)} &= -a_{vi}^{(k)} \sin \varphi^{(k)} + a_{vj}^{(k)} \cos \varphi^{(k)} \quad (v=1, 2, \dots, n). \end{aligned} \right\} \quad (3.6.8)$$

Аналогично, матрица $A^{(k+1)} = V_{ij}^{(k)'} B^{(k)}$ во всех строках, кроме i -й и j -й, будет иметь те же элементы, что и матрица $B^{(k)}$. Элементы i -й и j -й строк матрицы $A^{(k+1)}$ имеют соответственно вид

$$\left. \begin{aligned} a_{iv}^{(k+1)} &= b_{iv}^{(k)} \cos \varphi^{(k)} + b_{jv}^{(k)} \sin \varphi^{(k)}, \\ a_{jv}^{(k+1)} &= -b_{iv}^{(k)} \sin \varphi^{(k)} + b_{jv}^{(k)} \cos \varphi^{(k)} \quad (v=1, 2, \dots, n). \end{aligned} \right\} \quad (3.6.9)$$

Таким образом, для элемента $a_{ij}^{(k+1)}$ из формул (3.6.8) и (3.6.9) получим выражение

$$\begin{aligned} a_{ij}^{(k+1)} &= b_{ij}^{(k)} \cos \varphi^{(k)} + b_{jj}^{(k)} \sin \varphi^{(k)} = (-a_{ii}^{(k)} \sin \varphi^{(k)} + \\ &+ a_{ij}^{(k)} \cos \varphi^{(k)}) \cos \varphi^{(k)} + (-a_{ji}^{(k)} \sin \varphi^{(k)} + a_{jj}^{(k)} \cos \varphi^{(k)}) \sin \varphi^{(k)} = \\ &= a_{ij}^{(k)} \cos 2\varphi^{(k)} + \frac{1}{2} (a_{jj}^{(k)} - a_{ii}^{(k)}) \sin 2\varphi^{(k)}, \end{aligned} \quad (3.6.10)$$

так как $a_{ij}^{(k)} = a_{ji}^{(k)}$. Теперь из условия $a_{ij}^{(k+1)} = 0$ и формулы (3.6.10) определяем $\varphi^{(k)}$:

$$\operatorname{tg} 2\varphi^{(k)} = \frac{2a_{ij}^{(k)}}{a_{ii}^{(k)} - a_{jj}^{(k)}} \quad \left(|\varphi^{(k)}| \leq \frac{\pi}{4} \right) \quad (3.6.11)$$

или

$$\left. \begin{aligned} \cos \varphi^{(k)} &= \left[\frac{1}{2} \left(1 + (1 + p_k^2) \right)^{-\frac{1}{2}} \right]^{\frac{1}{2}}, \\ \sin \varphi^{(k)} &= \operatorname{sign} p_k \left[\frac{1}{2} \left(1 - (1 + p_k^2) \right)^{-\frac{1}{2}} \right]^{\frac{1}{2}}, \end{aligned} \right\} \quad (3.6.12)$$

$$\text{где } p_k = \frac{2a_{ij}^{(k)}}{a_{ii}^{(k)} - a_{jj}^{(k)}}.$$

Покажем теперь, что при таком выборе матрицы $V_{ij}^{(k)}$ максимально уменьшается сумма квадратов внедиагональных элементов матрицы $A^{(k+1)}$ по сравнению с соответствующей суммой матрицы $A^{(k)}$. Действительно, учитывая формулы (3.6.8) — (3.6.10) и симметричность матрицы $A^{(k)}$, получим

$$\begin{aligned} t^2(A^{(k+1)}) &= t^2(A^{(k)}) - 2[a_{ij}^{(k)}]^2 + \frac{1}{2} [(a_{jj}^{(k)} - a_{ii}^{(k)}) \sin 2\varphi^{(k)} + \\ &+ 2a_{ij}^{(k)} \cos 2\varphi^{(k)}]^2 = t^2(A^{(k)}) - 2[a_{ij}^{(k)}]^2 + \frac{1}{2} [2a_{ij}^{(k+1)}]^2 = \\ &= t^2(A^{(k)}) - 2[a_{ij}^{(k)}]^2, \end{aligned} \quad (3.6.13)$$

ибо, как уже отмечалось, $a_{ij}^{(k+1)} = 0$. Поскольку мы считаем, что $a_{ij}^{(k)}$ — максимальный по абсолютной величине внедиагональный элемент матрицы $A^{(k)}$, то из формулы (3.6.13) следует, что преобразование подобия с матрицей $V_{ij}^{(k)}$ максимально уменьшает величину $t^2(A^{(k)})$, т. е. является оптимальным для целей диагонализации матрицы $A^{(k)}$.

Теперь можно изложить алгоритм метода вращений и порядок вычислений:

1) в матрице $A^{(k)}$ ($k=0, 1, 2, \dots$) среди всех наддиагональных элементов выбираем максимальный по абсолютной величине элемент, определяем номера строки и столбца, в которых находится этот элемент, т. е. определяем числа i и j (если наибольших по абсолютной величине элементов несколько, то в качестве искомого можно взять любой или тот, для которого i имеет наименьшее значение);

2) по формуле (3.6.12) вычисляем $\cos \varphi^{(k)}$, $\sin \varphi^{(k)}$ и, используя названные числа i, j , находим по формулам (3.6.8), (3.6.9) элементы матрицы $A^{(k+1)}$;

3). итерационный процесс останавливаем, когда в пределах принятой точности величиной $t^2(A^{(k+1)})$ можно пренебречь, и в качестве приближенных собственных значений матрицы A с точностью до нумерации берем

$$\lambda_i = a_{ii}^{(k+1)} \quad (i=1, 2, \dots, n), \quad (3.6.14)$$

а в качестве собственных векторов — соответствующие столбцы матрицы

$$X^{(k+1)} = V_{i_0 j_0}^{(0)} V_{i_1 j_1}^{(1)} \dots V_{i_k j_k}^{(k)}. \quad (3.6.15)$$

3.6.2. Сходимость метода вращений

В этом пункте мы докажем сходимость метода вращений и получим некоторые сопутствующие этому вопросу оценки. При введенных нами обозначениях сходимость метода вращений означает, что

$$\lim_{k \rightarrow \infty} t^2(A^{(k)}) = 0.$$

Установим этот факт. В силу выбора элемента $a_{ij}^{(k)}$

$$t^2(A^{(k)}) \leq n(n-1) [a_{ij}^{(k)}]^2.$$

Значит,

$$[a_{ij}^{(k)}]^2 \geq \frac{t^2(A^{(k)})}{n(n-1)}. \quad (3.6.16)$$

Далее из формулы (3.6.13), используя неравенство (3.6.16), получим

$$t^2(A^{(k+1)}) = t^2(A^{(k)}) - 2[a_{ij}^{(k)}]^2 \leq t^2(A^{(k)}) - \frac{2t^2(A^{(k)})}{n(n-1)} = qt^2(A^{(k)}),$$

где

$$q = 1 - \frac{2}{n(n-1)} \quad \text{и} \quad 0 < q < 1,$$

ибо $n \geq 2$ (n — порядок матрицы A). Теперь легко выводятся такие неравенства:

$$\left. \begin{aligned} t^2(A^{(1)}) &\leq qt^2(A^{(0)}), \\ t^2(A^{(2)}) &\leq qt^2(A^{(1)}) \leq q^2 t^2(A^{(0)}), \\ &\dots \dots \dots \\ t^2(A^{(k)}) &\leq q^k t^2(A^{(0)}). \end{aligned} \right\} \quad (3.6.17)$$

Так как $t^2(A^{(0)}) = t^2(A) = \text{const} > 0$ и $0 < q < 1$, то из формулы (3.6.17) следует, что

$$\lim_{k \rightarrow \infty} t^2(A^{(k)}) = 0.$$

А это и означает, что описанный выше итерационный процесс сходится.

Перейдем к исследованию вопроса о скорости сходимости метода вращений и оценках погрешности в определении собственных значений и собственных векторов. С этой целью рассмотрим некоторую симметрическую матрицу A , внедиагональные элементы которой являются величинами не ниже первого порядка малости относительно $\varepsilon > 0$, и получим формулы, дающие возможность вычислять собственные значения и собственные векторы этой матрицы с точностью до некоторых величин, имеющих порядок малости относительно ε не ниже, чем второй.

Предположим, что названная матрица A — вещественная и ее собственные значения таковы, что внедиагональные элементы матрицы A малы по сравнению с числом ρ , определяемым по правилу

$$\rho = \min_{\lambda_i \neq \lambda_j} |\lambda_i - \lambda_j|; \quad (3.6.18)$$

где λ_i — собственные значения матрицы A . Очевидно, что если все λ_i равны между собой: $\lambda_1 = \lambda_2 = \dots = \lambda_n = \lambda$, то A — диагональная матрица, ибо в этом случае найдется такая ортогональная матрица V , что

$$V'AV = \Lambda = \lambda E,$$

и, следовательно,

$$A = \lambda EVV' = \lambda E,$$

ибо $VV' = E$ по определению ортогональной матрицы.

Пусть наибольшее по абсолютной величине собственное значение матрицы A имеет кратность m ($\lambda_1 = \lambda_2 = \dots = \lambda_m = \lambda$ и $|\lambda| > |\lambda_{m+1}| \geq \dots \geq |\lambda_n|$) и пусть диагональные элементы матрицы A расположены в порядке убывания их абсолютных величин. Разобьем матрицу A на клетки

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad (3.6.19)$$

таким образом, чтобы диагональные элементы A_{11} и только они были бы близкими к максимальному собственному значению λ . Порядок матрицы A_{11} должен, по-видимому, при этом совпасть с кратностью m собственного значения λ , если ε достаточно мало.

Для определения собственных векторов \bar{x} , соответствующих собственному значению λ , необходимо решить системы уравнений вида

$$(A - \lambda E)\bar{x} = \bar{0}. \quad (3.6.20)$$

Решения будут составлять подпространство размерности m , и вычислить их можно, например, следующим образом. Обозначим

$$\bar{x} = \begin{bmatrix} \bar{y} \\ \bar{z} \end{bmatrix},$$

где

$$y' = (x_1, x_2, \dots, x_m)', \quad z' = (x_{m+1}, x_{m+2}, \dots, x_n)'.$$

Систему (3.6.20) с учетом формулы (3.6.19) запишем в таком виде:

$$\left. \begin{aligned} (A_{11}-\lambda E)\bar{y}+A_{12}\bar{z} &= 0, \\ A_{21}\bar{y}+(A_{22}-\lambda E)\bar{z} &= 0. \end{aligned} \right\} \quad (3.6.21)$$

Первые m строк матрицы $A-\lambda E$ линейно зависимы с остальными, поэтому в качестве первого собственного вектора $\bar{x}_1 = \begin{bmatrix} \bar{y}_1 \\ \bar{z}_1 \end{bmatrix}$, отвечающего λ , можно взять вектор, для которого \bar{y}_1 определяется условно, как вектор $\underbrace{(1, 0, \dots, 0)'}_m$ а \bar{z}_1 — на основе формулы (3.6.21):

$$\bar{z}_1 = -(A_{22}-\lambda E)^{-1}A_{21}\bar{y}_1,$$

так как матрица $(A_{22}-\lambda E)^{-1}$ невырождена в силу формулы (3.6.18) и условий $|\lambda| > |\lambda_{m+1}| \geq \dots \geq |\lambda_n|$. Аналогично в качестве второго собственного вектора

$$\bar{x}_2 = \begin{bmatrix} \bar{y}_2 \\ \bar{z}_2 \end{bmatrix}, \text{ отвечающего } \lambda, \text{ можно взять вектор, для которого } \bar{y}_2 = \underbrace{(0, 1, 0, \dots, 0)'}_m,$$

а $\bar{z}_2 = -(A_{22}-\lambda E)^{-1}A_{21}\bar{y}_2$, и т. д. Продолжая этот процесс, мы в результате получим в качестве фундаментальной системы решений столбцы прямоугольной матрицы

$$X = \begin{bmatrix} E_m \\ -(A_{22}-\lambda E)^{-1}A_{21} \end{bmatrix} \quad (3.6.22)$$

размерности $n \times m$, где E_m — единичная матрица порядка m . Очевидно, что каждый из столбцов этой матрицы является собственным вектором матрицы A , отвечающим кратному собственному значению λ . Все столбцы матрицы X образуют в совокупности систему линейно независимых собственных векторов матрицы A , относящихся к λ .

Вычислим теперь в явном виде элементы x_{ij} ($m < i \leq n$, $j = 1, 2, \dots, m$) матрицы X . В силу формулы (3.6.22) и предположений относительно диагональных и внедиагональных элементов матрицы A имеем

$$x_{ij} = \frac{a_{ij}}{\lambda - a_{ii}} + O(\varepsilon^2), \quad (3.6.23)$$

где a_{ij} — элементы матрицы A_{21} .

Предположение об упорядоченности диагональных элементов матрицы A , которые мы сделали выше, можно опустить. Действительно, занумеруем собственные значения таким образом, чтобы λ_i было близко к диагональному элементу того же индекса, т. е. к a_{ii} . Аналогично тому, как это мы делали в предыдущем случае, можно показать, что собственные векторы матрицы A образуют соответственно столбцы матрицы

$$X = E + H, \quad (3.6.24)$$

где элементы h_{ij} матрицы H определяются по формулам

$$h_{ij} = \begin{cases} 0, & \text{если } \lambda_i = \lambda_j; \\ \frac{a_{ij}}{\lambda_j - a_{ii}} + O(\varepsilon^2), & \text{если } \lambda_i \neq \lambda_j. \end{cases} \quad (3.6.25)$$

Таким образом, если внедиагональные элементы вещественной симметрической матрицы имеют порядок малости не ниже первого относительно ε и малы по сравнению с числом ρ , то для собственных векторов такой матрицы имеют место формулы (3.6.24) и (3.6.25).

Выведем теперь аналогичные формулы для собственных значений λ_i . Обозначим через Λ диагональную матрицу собственных значений матрицы A .

Известно, что

$$X^{-1}AX = \Lambda.$$

Используя формулу (3.6.24), из этого равенства получим

$$\begin{aligned} A &= (E+H)\Lambda(E+H)^{-1} = (E+H)\Lambda(E-H+H^2-\dots) = \\ &= \Lambda + (H\Lambda - \Lambda H) - (H\Lambda - \Lambda H)H + O(H^3). \end{aligned} \quad (3.6.26)$$

Вычислим теперь элемент матрицы $\Lambda + (H\Lambda - \Lambda H) - (H\Lambda - \Lambda H)H$, стоящий на главной диагонали в i -й строке. Такими элементами будут: λ_i — у матрицы Λ , 0 — у матрицы

$(H\Lambda - \Lambda H)$ и $\sum_{k=1}^n h_{ki}h_{ik}(\lambda_k - \lambda_i)$ — у матрицы $(H\Lambda - \Lambda H)H$. Значит, искомый элемент есть сумма

$$\lambda_i + 0 - \sum_{k=1}^n h_{ki}h_{ik}(\lambda_k - \lambda_i),$$

и поэтому в силу формулы (3.6.26)

$$\lambda_i = a_{ii} + \sum_{k=1}^n h_{ki}h_{ik}(\lambda_k - \lambda_i)$$

или, учитывая формулы (3.6.25), окончательно получим

$$\lambda_i = a_{ii} + \sum_{p=1, p \in R_i}^n \frac{a_{ip}a_{pi}(\lambda_p - \lambda_i)}{(\lambda_p - a_{ii})(\lambda_i - a_{pp})} + O(\varepsilon^3), \quad (3.6.27)$$

где R_i — множество тех чисел j из ряда $1, 2, \dots, n$, для которых $\lambda_j = \lambda_i$.

В этой формуле отношения

$$\frac{\lambda_p - \lambda_i}{(\lambda_p - a_{ii})(\lambda_i - a_{pp})}$$

ограничены сверху, ибо при достаточно малом ε $|\lambda_p - a_{ii}| \geq \rho$, $|\lambda_i - a_{pp}| \geq \rho$ и $\rho > 0$. Так как внедиагональные элементы матрицы A по абсолютной величине имеют порядок малости не ниже первого относительно ε , то из формулы (3.6.27) теперь следует, что

$$\lambda_i = a_{ii} + O(\varepsilon^2), \quad (3.6.28)$$

независимо от кратности корня λ_i . Этим выражением можно воспользоваться при записи формул (3.6.25) и (3.6.27), что дает такое правило для вычисления элементов h_{ij} матрицы H и собственных чисел матрицы A :

$$h_{ij} = \begin{cases} 0, & \text{если } \lambda_i \neq \lambda_j; \\ \frac{a_{ij}}{a_{jj} - a_{ii}} + O(\varepsilon^2), & \text{если } \lambda_i = \lambda_j \end{cases} \quad (3.6.29)$$

и

$$\lambda_i = a_{ii} + \sum_{p=1}^n \frac{a_{ip}a_{pi}}{a_{ii} - a_{pp}} + O(\epsilon^3). \quad (3.6.30)$$

Все вышеизложенное дает нам основание сформулировать следующую теорему, доказанную впервые В. В. Воеводиным в более общем виде, чем тот, который мы здесь приводим.

Теорема 1. Если матрица A вещественная и симметрическая и ее внедиагональные элементы являются величинами не ниже первого порядка малости относительно ϵ ,*) то для матрицы X , составленной из собственных векторов матрицы A , и собственных значений λ_i справедливы формулы (3.6.24), (3.6.29) и (3.6.30).

Укажем здесь еще на одно важное обстоятельство, а именно, на особое свойство элементов a_{ij} матрицы A , для которых $\lambda_i = \lambda_j$. Оказывается, что такие элементы матрицы A будут величинами не ниже второго порядка малости относительно ϵ . Действительно, будем вычислять внедиагональные элементы матрицы A по формуле (3.6.26), считая известными матрицы Λ и H . Так как $h_{ij} = 0$ при $\lambda_i = \lambda_j$ и $h_{ij} = \frac{a_{ij}}{a_{jj} - a_{ii}} + O(\epsilon^2)$ при $\lambda_i \neq \lambda_j$, то из формулы (3.6.26) следует, что a_{ij} есть величина не ниже второго порядка малости относительно ϵ при $\lambda_i = \lambda_j$.

Покажем теперь, как эти результаты могут быть применены к исследованию скорости сходимости метода вращений в случае вещественной симметрической матрицы A .

Пусть $A^{(k)}$ есть матрица, полученная после выполнения k -го шага метода вращений, и пусть $a_{ij}^{(k)}$ — ее элементы. Предположим, что все внедиагональные элементы этой матрицы являются величинами не ниже первого порядка малости относительно некоторого достаточно малого числа $\epsilon > 0$. Тогда с учетом симметричности матрицы $A^{(k)}$ в силу формулы (3.6.30) получим

$$\lambda_i = a_{ii}^{(k)} + \sum_{p=1, p \in R_i} \frac{[a_{ip}^{(k)}]^2}{a_{ii}^{(k)} - a_{pp}^{(k)}} + O(\epsilon^3).$$

Матрицу V собственных векторов исходной матрицы A можно вычислить таким путем. По смыслу алгоритма метода вращений имеем

$$A^{(k)} = T'AT,$$

где $T = V_{i_0 j_0}^{(0)} V_{i_1 j_1}^{(1)} \dots V_{i_k j_k}^{(k)}$. Применительно к матрице $A^{(k)}$ можно вычислить такую матрицу X по формулам (3.6.24) и (3.6.29), заменяя в последней формуле a_{ij} на $a_{ij}^{(k)}$, что будет справедливо равенство

$$X^{-1}A^{(k)}X = \Lambda,$$

где Λ — диагональная матрица собственных значений искомой матрицы A . Значит,

$$X^{-1}A^{(k)}X = X^{-1}T^{-1}ATX = \Lambda, \quad (3.6.31)$$

*) В качестве величины, характеризующей малость внедиагональных элементов неособенной матрицы, можно, например, взять отношение $\epsilon = \frac{M}{m}$, где $M = \max_{i \neq j} |a_{ij}|$, $m = \min_i |a_{ii}|$. При преобразованиях вращения эта величина будет меняться от шага к шагу и при увеличении числа преобразований будет стремиться к нулю.

ибо T — ортогональная матрица и для нее $T' = T^{-1}$. Поэтому из (3.6.31) следует, что искомая матрица V равна

$$V = TX = V_{i_0 j_0}^{(0)} V_{i_1 j_1}^{(1)} \dots V_{i_k j_k}^{(k)} X.$$

Докажем теперь теорему о скорости сходимости метода вращений.

Теорема 2. *Какова бы ни была вещественная симметрическая матрица A , метод вращений для нее обладает квадратичной сходимостью.*

Доказательство. Предположим, что процесс метода вращений проведен настолько далеко, что все внедиагональные элементы матрицы $A^{(k)}$ стали величинами не ниже первого порядка малости относительно ε и малыми по сравнению с величиной ρ , определяемой по формуле (3.6.18). Такое число k существует, ибо мы показали ранее,

что $t^2(A^{(k)}) \rightarrow 0$ при $k \rightarrow \infty$. Мы уже отмечали, что элементы $a_{ij}^{(k)}$, соответствующие собственным значениям $\lambda_i = \lambda_j$, являются величинами не ниже второго порядка малости относительно ε . Поэтому среди этих элементов максимального внедиагонального элемента матрицы $A^{(k)}$ не будет, ибо по идее метода вращений такой элемент на рассматриваемом шаге преобразований должен быть величиной первого порядка малости относительно ε . Мы предположим, что максимальный элемент, например $a_{ij}^{(k)}$ ($i < j$), будет величиной порядка ε . Чтобы в методе вращений аннулировать этот элемент, надо для соответствующей матрицы вращения определить $\cos \varphi^{(k)}$ и $\sin \varphi^{(k)}$. Из формул (3.6.12) следует, что в этом случае

$$\cos \varphi^{(k)} = 1 + O(\varepsilon^2), \quad \sin \varphi^{(k)} = O(\varepsilon),$$

ибо

$$\rho_k = \frac{2a_{ij}^{(k)}}{a_{ii}^{(k)} - a_{jj}^{(k)}} = \frac{2a_{ij}^{(k)}}{\lambda_i - \lambda_j + O(\varepsilon^2)}$$

и число $\lambda_i - \lambda_j + O(\varepsilon^2)$ ограничено снизу, а $a_{ij}^{(k)}$ имеет по предположению первый порядок относительно ε . Ясно, что с таким образом подобранным углом поворота $\varphi^{(k)}$ все элементы, меняющиеся на $(k+1)$ -м шаге преобразований подобия могут изменяться (см. формулы (3.6.8) и (3.6.9)) только на величины порядка ε^2 . При этом $a_{ij}^{(k+1)}$ и $a_{ji}^{(k+1)}$ равны

нулю по построению матрицы вращения. Заметим, что элемент матрицы $A^{(k)}$, аннулированный на предыдущем шаге, может стать не более величины порядка ε^2 . Учитывая то, что A симметричная матрица и симметричны все матрицы вида $A^{(k)}$, мы сможем не более

чем за r_1 $\left(r_1 \leq \frac{n(n+1)}{2} \right)$ элементарных шагов все внедиагональные элементы $\left(\text{число их равно } \frac{n(n-1)}{2} \right)$ сделать величинами не ниже второго порядка малости относительно ε . Таким образом, если $t^2(A^{(k)}) \leq O(\varepsilon)$, то $t^2(A^{(k_1)}) \leq O(\varepsilon^2)$, где $k_1 = k + r_1$. Продолжая по аналогии этот процесс, получим

$$t^2(A^{(k_m)}) \leq O(\varepsilon^{2^m}),$$

где

$$k_m = k + r_1 + r_2 + \dots + r_m \quad \text{и} \quad r_i \leq \frac{n(n-1)}{2} \quad (1 \leq i \leq m).$$

Это и доказывает квадратичную сходимость метода вращений.

Другие строки матрицы $A^{(k+1)}$ будут иметь те же элементы, что и строки матрицы $B^{(k)}$.

Так как матрица A эрмитова, то $a_{ij} = \overline{a_{ji}}$, а следовательно, и $a_{ij}^{(k)} = \overline{a_{ji}^{(k)}}$. Поэтому на основании формул (3.6.33) и (3.6.34) между величинами $t^2(A^{(k+1)})$ и $t^2(A^{(k)})$ может быть установлена следующая связь:

$$t^2(A^{(k+1)}) = t^2(A^{(k)}) - 2|a_{ij}^{(k)}|^2 + 2 \left[|a_{ij}^{(k)}| (e^{i\alpha^{(k)}} \cos^2 \varphi^{(k)} - e^{i(2\psi^{(k)} - \alpha^{(k)})} \sin^2 \varphi^{(k)}) + (a_{jj}^{(k)} - a_{ii}^{(k)}) e^{i\psi^{(k)}} \frac{1}{2} \sin 2\varphi^{(k)} \right]^2, \quad (3.6.35)$$

где обозначено $\alpha^{(k)} = \arg a_{ij}^{(k)}$. В этой формуле $a_{ij}^{(k)}$, как уже отмечалось, — наибольший по модулю внедиагональный элемент матрицы $A^{(k)}$. Значит, из формулы (3.6.35) следует, что для максимального уменьшения $t^2(A^{(k+1)})$ параметры $\varphi^{(k)}$ и $\psi^{(k)}$ следует выбирать так, чтобы выполнялось условие

$$\left| |a_{ij}^{(k)}| (e^{i\alpha^{(k)}} \cos^2 \varphi^{(k)} - e^{i(2\psi^{(k)} - \alpha^{(k)})} \sin^2 \varphi^{(k)}) + (a_{jj}^{(k)} - a_{ii}^{(k)}) e^{i\psi^{(k)}} \frac{1}{2} \sin 2\varphi^{(k)} \right| = 0.$$

Это условие будет выполнено, если мы положим

$$\psi^{(k)} = \arg a_{ij}^{(k)}, \quad \operatorname{tg} 2\varphi^{(k)} = \frac{2|a_{ij}^{(k)}|}{a_{ii}^{(k)} - a_{jj}^{(k)}} \quad \left(|\varphi^{(k)}| \leq \frac{\pi}{4} \right).$$

Определив таким образом $\psi^{(k)}$ и $\varphi^{(k)}$ и построив с этими значениями матрицу $V_{ij}(\varphi^{(k)}, \psi^{(k)})$, мы сможем затем выполнить элементарный шаг преобразований подобия матрицы $A^{(k)}$. Это преобразование для целей диагонализации матрицы $A^{(k)}$ будет оптимальным, ибо оно максимально уменьшает $t^2(A^{(k)})$.

Полученные расчетные формулы несущественно отличаются от аналогичных формул в случае вещественных симметрических матриц. Отметим, что все результаты, которые мы получили по сходимости и скорости сходимости метода вращений, могут быть перенесены и на эрмитовы матрицы [2].

Если матрица косоэрмитова, т. е. матрица, удовлетворяющая условию $A = -A^*$, то к ней метод вращений можно применить после несложной замены. А именно, вместо матрицы A рассмотрим матрицу B , определяемую по правилу

$$B = iA.$$

Легко проверить, что $B=B^*$, и, значит, B — эрмитова матрица. К ней можно применить метод вращений и найти таким образом ее собственные значения $\lambda_k(B)$. Но $\lambda_k(A) = -i\lambda_k(B)$, чем и заканчивается задача о вычислении собственных значений $\lambda_k(A)$ косоэрмитовой матрицы A .

В заключение укажем на некоторые особенности метода вращений при его реализации на ЭВМ. Вычислительная практика показала высокую надежность и нечувствительность метода вращений к таким свойствам матрицы, как расположение по абсолютной величине собственных значений, в частности их близость и кратность. Это делает метод вращений одним из наиболее эффективных методов решения полной проблемы собственных значений.

Недостатком описанного здесь алгоритма является то, что на каждом шаге приходится отыскивать максимальный по модулю внедиагональный элемент. На эту операцию при использовании ЭВМ затрачивается много машинного времени, что снижает эффективность алгоритма.

Поэтому заслуживает внимания следующий выбор внедиагонального элемента, подлежащего аннулированию на $(k+1)$ -м шаге.

Как мы видели, если аннулируемым на $(k+1)$ -м шаге был элемент $a_{ij}^{(k)}$, то в силу формулы (3.6.35) должны иметь место равенства

$$\sum_{m=1, m \neq p}^n |a_{pm}^{(k+1)}|^2 = \sum_{m=1, m \neq p}^n |a_{pm}^{(k)}|^2$$

для всех $p \neq i, j$. Эти равенства говорят о том, что если в начале процесса вычислений по матрице $A^{(0)}=A$ составить суммы квадратов модулей внедиагональных элементов каждой строки и обозначить их через

$$\sigma_1, \sigma_2, \dots, \sigma_n \quad (\sigma_p = \sum_{m=1, m \neq p}^n |a_{pm}|^2, (p=1, 2, \dots, n)),$$

то при выполнении в дальнейшем над матрицей $A^{(0)}=A$ элементарного шага преобразований с матрицей $V_{ij}^{(0)}(\varphi^{(0)}, \psi^{(0)})$ будут меняться только два числа, а именно σ_i и σ_j , другие же числа σ_p останутся без изменения.

Этот факт позволяет находить почти максимальный (будем называть его оптимальным) элемент матрицы A путем просмотра всего лишь $2n-1$ чисел. Делается это так. Сначала находим суммы $\sigma_1, \sigma_2, \dots, \sigma_n$ и выбираем среди них наибольшую. Пусть это будет, например, σ_i . Ясно, что для нахождения этой суммы надо выполнить n просмотров. Далее по индексу этой суммы находим i -ю строку в матрице A и в этой строке находим наибольший по модулю элемент. Это можно сделать, выполнив $n-1$ просмотров. Пусть таким элементом будет a_{ij} ($j \neq i$). Его и принимают за искомым оптимальным элементом, подлежащий аннулированию.

Конечно, из сказанного здесь не следует, что оптимальный элемент будет обязательно наибольшим по модулю внедиагональным элементом матрицы A . Однако он будет близким к такому элементу и во всяком случае будет не менее среднего квадратического всех внедиагональных элементов. После аннулирования оптимального элемента a_{ij} суммы $\sigma_1, \sigma_2, \dots, \sigma_n$ подготавливаются к следующему шагу путем пересчета только σ_i и σ_j , другие суммы остаются прежними.

Отметим, что теория метода вращений с выбором максимального элемента переносится на вариант метода вращений с выбором оптимального элемента [2].

§ 3.7. УТОЧНЕНИЕ СОБСТВЕННЫХ ЗНАЧЕНИЙ И ПРИНАДЛЕЖАЩИХ ИМ СОБСТВЕННЫХ ВЕКТОРОВ МАТРИЦ И УСКОРЕНИЕ СХОДИМОСТИ МЕТОДА ИТЕРАЦИИ ПРИ РЕШЕНИИ СИСТЕМ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ

При решении полной и частичной проблем собственных значений иногда возникает необходимость уточнить полученные результаты, и в таких случаях важно располагать приемами, позволяющими уточнять полученные приближенные собственные значения и отвечающие им собственные векторы. Необходимость уточнения может быть обусловлена различными причинами, например при решении полной проблемы собственных значений методами, основанными на получении характеристического полинома матрицы, может оказаться, что из-за неустойчивости метода к ошибкам округления такой многочлен будет получен с неудовлетворительной точностью и, следовательно, нельзя будет хорошо вычислить и собственные значения матрицы.

Аналогично если решается частичная проблема, то, как мы видели, большинство рассматриваемых методов в этом случае имеет медленную сходимость, не превышающую скорости сходимости геометрической прогрессии. Это обстоятельство заставляет выполнять большое число итераций, для того чтобы получить результат с требуемой точностью. Здесь возникает проблема ускорения сходимости медленно сходящихся последовательностей. Такая же проблема возникает и для некоторых итерационных методов, предназначенных для решения систем линейных алгебраических уравнений.

Изучению этих вопросов и посвящается настоящий параграф.

3.7.1. Уточнение полной проблемы собственных значений

Пусть матрица A имеет попарно различные собственные значения, которые удовлетворяют, например, таким неравенствам:

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|.$$

Предположим, что, используя какой-либо из методов, рассмотренных нами в этой главе, мы нашли приближенные собственные значения

$\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_n$ матрицы A и отвечающие им приближенные собственные векторы $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$. Пусть нам известны также и приближенные собственные векторы $\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n$ сопряженной матрицы A^* .

Поставим задачу уточнить всю совокупность перечисленных приближенных величин. Уточненные значения λ_i, \bar{x}_i и \bar{y}_i будем разыскивать в виде

$$\left. \begin{aligned} \lambda_i &= \tilde{\lambda}_i + \Delta\tilde{\lambda}_i, \\ \bar{x}_i &= \tilde{x}_i + \Delta\tilde{x}_i, \\ \bar{y}_i &= \tilde{y}_i + \Delta\tilde{y}_i \\ (i=1, 2, \dots, n), \end{aligned} \right\} \quad (3.7.1)$$

где величины $\Delta\tilde{\lambda}_i, \Delta\tilde{x}_i, \Delta\tilde{y}_i$ имеют смысл поправок, которые мы будем считать малыми и которые нам предстоит определить.

Запишем разложение векторов $\tilde{x}_i + \Delta\tilde{x}_i$ и $\tilde{y}_i + \Delta\tilde{y}_i$ соответственно по базисам $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$ и $\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n$:

$$\left. \begin{aligned} \tilde{x}_i + \Delta\tilde{x}_i &= \sum_{j=1}^n h_{ij} \tilde{x}_j, \\ \tilde{y}_i + \Delta\tilde{y}_i &= \sum_{j=1}^n g_{ij} \tilde{y}_j. \end{aligned} \right\} \quad (3.7.2)$$

Поскольку собственные векторы $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$ и $\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n$ определяются с точностью до произвольного множителя, то можно, не нарушая общности, положить в приведенных выше разложениях коэффициенты h_{ii} и g_{ii} равными единице. Тогда из формулы (3.7.2) получим

$$\left. \begin{aligned} \Delta\tilde{x}_i &= \sum_{j=1, j \neq i}^n h_{ij} \tilde{x}_j, \\ \Delta\tilde{y}_i &= \sum_{j=1, j \neq i}^n g_{ij} \tilde{y}_j. \end{aligned} \right\} \quad (3.7.3)$$

Так как $\Delta\tilde{x}_i$ и $\Delta\tilde{y}_i$ по предположению малы, то в силу формулы (3.7.3) малыми, вообще говоря, должны быть и числа h_{ij}, g_{ij} . Эта формула показывает, что если мы сможем определить числа h_{ij} и g_{ij} , то тем самым будут определены и искомые поправки $\Delta\tilde{x}_i, \Delta\tilde{y}_i$. Поэтому сейчас перед нами встает проблема нахождения коэффициентов h_{ij}, g_{ij} и поправки $\Delta\tilde{\lambda}_i$.

Укажем способ вычисления $\Delta\tilde{\lambda}_i$ и коэффициентов h_{ij} и g_{ij} через невязки известного нам приближенного решения полной проблемы собственных значений, т. е. через величины

$$\left. \begin{aligned} \bar{r}_i &= A\tilde{x}_i - \tilde{\lambda}_i \tilde{x}_i, \\ \bar{s}_i &= A^* \tilde{y}_i - \tilde{\lambda}_i \tilde{y}_i, \end{aligned} \right\} \quad (3.7.4)$$

где $\tilde{\lambda}_i$ — числа, комплексно сопряженные с $\tilde{\lambda}_i$.

Будем находить поправки $\Delta\tilde{\lambda}_i$, $\Delta\tilde{x}_i$, $\Delta\tilde{y}_i$ таким образом, чтобы выполнялись равенства

$$A\bar{x}_i = \lambda_i \bar{x}_i, \quad (3.7.5)$$

$$A^* \bar{y}_i = \bar{\lambda}_i \bar{y}_i, \quad (3.7.6)$$

где λ_i , \bar{x}_i , \bar{y}_i определяются по формулам (3.7.1).

Перепишем формулу (3.7.5) в таком виде:

$$A(\tilde{x}_i + \Delta\tilde{x}_i) = (\tilde{\lambda}_i + \Delta\tilde{\lambda}_i)(\tilde{x}_i + \Delta\tilde{x}_i).$$

Это соотношение можно рассматривать как нелинейное уравнение относительно неизвестных $\Delta\tilde{x}_i$ и $\Delta\tilde{\lambda}_i$. Удерживая в этом выражении только линейные члены относительно $\Delta\tilde{x}_i$ и $\Delta\tilde{\lambda}_i$ и используя формулу (3.7.4), получим из него линейное уравнение

$$A\Delta\tilde{x}_i - \tilde{\lambda}_i \Delta\tilde{x}_i \approx -r_i + \Delta\tilde{\lambda}_i \tilde{x}_i. \quad (3.7.7)$$

Отметим, что если $\Delta\tilde{x}_i$ и $\Delta\tilde{\lambda}_i$ достаточно малы, то уравнение (3.7.7) позволит определить главные части точных значений погрешностей, и отбрасывание членов второго порядка малости, которое мы совершили, не повлияет сильно на окончательный результат. Образует теперь с помощью формулы (3.7.7) и векторов \tilde{y}_j ($j=1, 2, \dots, n$) скалярные произведения

$$(A\Delta\tilde{x}_i, \tilde{y}_j) - \tilde{\lambda}_i(\Delta\tilde{x}_i, \tilde{y}_j) = -(\bar{r}_i, \tilde{y}_j) + \Delta\tilde{\lambda}_i(\tilde{x}_i, \tilde{y}_j). \quad (3.7.8)$$

Из уравнения (3.7.8) неизвестное $\tilde{\Delta x}_i$ можно исключить следующим путем. По определению

$$(A\tilde{\Delta x}_i, \tilde{y}_j) = (\tilde{\Delta x}_i, A^* \tilde{y}_j) = (\tilde{\Delta x}_i, A^*(\bar{y}_j - \tilde{\Delta y}_j)).$$

Значит, с точностью до величин второго порядка малости верно равенство

$$(A\tilde{\Delta x}_i, \tilde{y}_j) - \tilde{\lambda}_j(\tilde{\Delta x}_i, \tilde{y}_j) \approx 0. \quad (3.7.9)$$

Положим в формулах (3.7.8) и (3.7.9) $j=i$, тогда для $\tilde{\Delta \lambda}_i$ можно получить такое выражение:

$$\tilde{\Delta \lambda}_i \approx \frac{(\bar{r}_i, \tilde{y}_i)}{(\tilde{x}_i, \tilde{y}_i)}. \quad (3.7.10)$$

Из этой формулы при $i=1, 2, \dots, n$ вычислим искомые поправки $\tilde{\Delta \lambda}_1, \tilde{\Delta \lambda}_2, \dots, \tilde{\Delta \lambda}_n$, после чего уточненные значения собственных чисел $\lambda_1, \lambda_2, \dots, \lambda_n$ найдем по формуле (3.7.1).

Укажем теперь правила для вычисления поправок $\tilde{\Delta x}_i$ и $\tilde{\Delta y}_i$. В силу формулы (3.7.3) для этого достаточно указать правила для вычисления коэффициентов h_{ij} и g_{ij} ($j=1, 2, \dots, n$). Пусть $i \neq j$, тогда, используя формулу (3.7.3), получим

$$(\tilde{\Delta x}_i, \tilde{y}_j) \approx h_{ij}(\tilde{x}_j, \tilde{y}_j) \quad \text{и} \quad \tilde{\Delta \lambda}_i(\tilde{x}_i, \tilde{y}_j) \approx 0. \quad (3.7.11)$$

Эти приближенные равенства выводятся следующим путем. По определению величин $\tilde{\Delta x}_i, \tilde{x}_i$ и \tilde{y}_i имеем

$$\begin{aligned} (\tilde{\Delta x}_i, \tilde{y}_j) &= \left(\sum_{s=1, s \neq i}^n h_{is} \tilde{x}_s, \tilde{y}_j \right) = \sum_{s=1, s \neq i}^n h_{is} (\bar{x}_s - \tilde{\Delta x}_s, \bar{y}_j - \tilde{\Delta y}_j) = \\ &= \sum_{s=1, s \neq i}^n h_{is} (\bar{x}_s, \bar{y}_j) + \sum_{s=1, s \neq i}^n h_{is} [-(\tilde{\Delta x}_s, \bar{y}_j) - (\bar{x}_s, \tilde{\Delta y}_j) + (\tilde{\Delta x}_s, \tilde{\Delta y}_j)]. \end{aligned}$$

Учитывая, что $(\bar{x}_s, \bar{y}_j) = 0$ при $s \neq j$, для скалярного произведения $(\tilde{\Delta x}_i, \tilde{y}_j)$ получим:

$$\begin{aligned}
(\Delta \tilde{x}_i, \tilde{y}_j) &= h_{ij}(\bar{x}_j, \bar{y}_j) + h_{ij} [-(\Delta \tilde{x}_j, \bar{y}_j) - (\bar{x}_j, \Delta \tilde{y}_j) + (\Delta \tilde{x}_j, \Delta \tilde{y}_j)] + \\
&+ \sum_{s=1, s \neq i, j}^n h_{is} [-(\Delta \tilde{x}_s, \bar{y}_j) - (\bar{x}_s, \Delta \tilde{y}_j) + (\Delta \tilde{x}_s, \Delta \tilde{y}_j)] = h_{ij}(\tilde{x}_j, \tilde{y}_j) + \\
&+ \sum_{s=1, s \neq i, j}^n h_{is} [-(\Delta \tilde{x}_s, \bar{y}_j) - (\bar{x}_s, \Delta \tilde{y}_j) + (\Delta \tilde{x}_s, \Delta \tilde{y}_j)].
\end{aligned}$$

Отсюда следует первое приближенное равенство из формулы (3.7.11). Второе приближенное равенство становится очевидным, если учесть условие $(\bar{x}_i, \bar{y}_j) = 0$ при $i \neq j$ и формулы (3.7.1).

Далее с помощью формулы (3.7.9) из уравнения (3.7.8) получим

$$\tilde{\lambda}_j(\Delta \tilde{x}_i, \tilde{y}_j) - \tilde{\lambda}_i(\Delta \tilde{x}_i, \tilde{y}_j) \approx -(\bar{r}_i, \tilde{y}_j) + \Delta \tilde{\lambda}_i(\tilde{x}_i, \tilde{y}_j).$$

Отбрасывая здесь величину $\Delta \tilde{\lambda}_i(\tilde{x}_i, \tilde{y}_j)$, имеющую второй порядок малости относительно $\Delta \tilde{\lambda}_i, \Delta \tilde{x}_i, \Delta \tilde{y}_j$, и учитывая, что $(\Delta \tilde{x}_i, \tilde{y}_j) \approx h_{ij}(\tilde{x}_j, \tilde{y}_j)$, окончательно получим следующее приближенное уравнение для определения h_{ij} :

$$(\tilde{\lambda}_j - \tilde{\lambda}_i) h_{ij}(\tilde{x}_j, \tilde{y}_j) \approx -(\bar{r}_i, \tilde{y}_j).$$

Таким образом, коэффициент h_{ij} можно находить по формуле

$$h_{ij} \approx \frac{(\bar{r}_i, \tilde{y}_j)}{(\tilde{\lambda}_i - \tilde{\lambda}_j)(\tilde{x}_j, \tilde{y}_j)} \quad (j=1, 2, \dots, n, j \neq i). \quad (3.7.12)$$

Аналогично из формулы (3.7.6) можно получить правило также и для вычисления коэффициента g_{ij} :

$$g_{ij} \approx \frac{(\bar{s}_i, \tilde{x}_j)}{(\tilde{\lambda}_i - \tilde{\lambda}_j)(\tilde{x}_j, \tilde{y}_j)} \quad (j=1, 2, \dots, n, j \neq i). \quad (3.7.13)$$

После того как будут найдены коэффициенты h_{ij} и g_{ij} , по формуле (3.7.3) определяются поправки $\Delta \tilde{x}_i$ и $\Delta \tilde{y}_i$ и по правилу (3.7.1) — уточненные значения соответствующих собственных векторов \tilde{x}_i и \tilde{y}_i . Заметим, что для определения поправок $\Delta \tilde{\lambda}_i$ и коэффициентов h_{ij}, g_{ij} , кроме

исходных величин $\tilde{\lambda}_i, \tilde{x}_i, \tilde{y}_i$, необходимо знать еще невязки \bar{r}_i и \bar{s}_i , т. е. те величины, которые получаются в результате контрольных вычислений.

Если при контрольных вычислениях окажется, что невязка \bar{r}_i велика, то вторую невязку \bar{s}_i можно по формуле (3.7.4) не вычислять, ибо она нужна только для вычисления коэффициентов g_{ij} . А эти коэффициенты после определения h_{ij} можно находить, как в этом легко убедиться, по формуле

$$\bar{g}_{ij} = - \frac{(\tilde{x}_i, \tilde{y}_j) + h_{ij}(\tilde{x}_j, \tilde{y}_j)}{(\tilde{x}_i, \tilde{y}_i)},$$

где \bar{g}_{ij} — число, комплексно сопряженное с g_{ij} .

Избранный нами процесс определения поправок $\Delta\tilde{\lambda}_i, \Delta\tilde{x}_i$ и $\Delta\tilde{y}_i$, как можно убедиться, эквивалентен применению одного шага метода Ньютона (гл. 1) к нелинейной системе уравнений (3.7.5) и (3.7.6).

Если матрица A — вещественная и симметрическая, то $\tilde{x}_i = \tilde{y}_i$, и формулы (3.7.10) и (3.7.12) упрощаются:

$$\begin{aligned}\Delta\tilde{\lambda}_i &\approx \frac{(\bar{r}_i, \tilde{x}_i)}{(\tilde{x}_i, \tilde{x}_i)}, \\ h_{ij} &\approx \frac{(\bar{r}_i, \tilde{x}_j)}{(\tilde{\lambda}_i - \tilde{\lambda}_j)(\tilde{x}_j, \tilde{x}_j)}.\end{aligned}$$

Изложенный здесь способ уточнения собственных значений и векторов матрицы A следует повторить, если полученные уточненные величины будут неудовлетворительными по точности.

3.7.2. Уточнение отдельного собственного значения и принадлежащего ему собственного вектора

Пусть мы нашли приближенное собственное значение $\tilde{\lambda}$ и отвечающий ему собственный вектор \tilde{x} некоторой матрицы A . Соответствующие точные значения обозначим через λ и \bar{x} . Пусть

$$\begin{aligned}\lambda &= \tilde{\lambda} + \Delta\tilde{\lambda}, \\ \bar{x} &= \tilde{x} + \Delta\tilde{x}.\end{aligned}$$

Матрица R существует. Действительно, если потребовать, чтобы условие

$$R(A_0 - \lambda_0 E)\bar{y} = \bar{y}$$

выполнялось на подпространстве векторов \bar{y} , ортогональных к вектору \bar{x}_0 , то это будет эквивалентно образованию системы $n^2 - n$ линейных уравнений для определения n^2 элементов матрицы R :

$$R(A_0 - \lambda_0 E)\bar{y}_j = \bar{y}_j \quad (j = 1, 2, \dots, n-1),$$

где $(\bar{y}_j, \bar{x}_0) = 0$ и векторы $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_{n-1}$ линейно независимы. Если учесть еще условие $R\bar{x}_0 = 0$, то для определения элементов матрицы R мы получим линейную систему n^2 уравнений вида

$$B_0 \bar{r} = \bar{b},$$

где приняты обозначения:

$$B_0 = \begin{bmatrix} z & & & 0 \\ & z & & \\ & & \ddots & \\ 0 & & & z \end{bmatrix}, \quad z = \begin{bmatrix} z_{11} & z_{21} & \dots & z_{n1} \\ z_{12} & z_{22} & \dots & z_{n2} \\ \dots & \dots & \dots & \dots \\ z_{1n-1} & z_{2n-1} & \dots & z_{nn-1} \\ x_{10} & x_{20} & \dots & x_{n0} \end{bmatrix},$$

$$\bar{z}_j = (A_0 - \lambda_0 E)\bar{y}_j, \quad \bar{y}_j = (y_{1j}, y_{2j}, \dots, y_{nj})',$$

$$\bar{z}_j = (z_{1j}, z_{2j}, \dots, z_{nj})', \quad \bar{x}_0 = (x_{10}, x_{20}, \dots, x_{n0})',$$

$$\bar{r} = (r_{11}, r_{12}, \dots, r_{1n}, r_{21}, r_{22}, \dots, r_{2n}, \dots, r_{n1}, r_{n2}, \dots, r_{nn})',$$

$$\bar{b} = (y_{11}, y_{12}, \dots, y_{1n-1}, 0, y_{21}, y_{22}, \dots, y_{2n-1}, 0, \dots, y_{n1}, y_{n2}, \dots, y_{nn-1}, 0)'. \quad .$$

Матрица $A_0 - \lambda_0 E$ неособенна на подпространстве векторов \bar{y}_j , ортогональных к вектору \bar{x}_0 , т. е. $(A_0 - \lambda_0 E)\bar{y}_j \neq 0$. В самом деле, если бы $(A_0 - \lambda_0 E)\bar{y}_j = 0$, то тогда было бы $A_0 \bar{y}_j = \lambda_0 \bar{y}_j$ и \bar{y}_j являлся бы собственным вектором матрицы A_0 , отвечающим λ_0 . А это невозможно, так как $(\bar{y}_j, \bar{x}_0) = 0$ и λ_0 является простым собственным значением матрицы A_0 .

Векторы \bar{y}_j и \bar{x}_0 линейно независимы, в силу вышеизложенного линейно независимыми будут и векторы \bar{z}_j, \bar{x}_0 . Поэтому $\det z \neq 0$ и матрица B_0 неособенна.

Заметим, что мы только выяснили здесь вопрос существования матрицы R . Как мы увидим далее, в рассматриваемом правиле улучшения собственного значения не будет необходимости находить элементы матрицы R из системы $B_0 \bar{r} = \bar{b}$.

Будем считать также, что λ — собственное значение матрицы A , а \bar{x} — соответствующий ему собственный вектор и что для векторов \bar{x}_0 и \bar{x} выполняется условие

$$(\bar{x}, \bar{x}_0) \neq 0.$$

При этом можно полагать, что $\bar{x} = \bar{x}_0 + z$ и $(z, \bar{x}_0) = 0$.

Введем в рассмотрение матрицу B и число μ , определяемые по формулам

$$B = A - A_0,$$

$$\mu = \lambda - \lambda_0,$$

и поставим задачу найти вектор \bar{x} и число μ , используя величины R , B и \bar{x}_0 .

Сначала проверим справедливость следующих равенств:

$$((\mu E - B)\bar{x}, \bar{x}_0) = 0 \quad (3.7.16)$$

и

$$(E - R(\mu E - B))\bar{x} = \bar{x}_0. \quad (3.7.17)$$

Имеем

$$(\mu E - B)\bar{x} = (\lambda E - A)\bar{x} - (\lambda_0 E - A_0)\bar{x} = (A_0 - \lambda_0 E)\bar{x},$$

ибо $(\lambda E - A)\bar{x} = 0$ по определению λ и \bar{x} . Далее получим

$$((\mu E - B)\bar{x}, \bar{x}_0) = ((A_0 - \lambda_0 E)\bar{x}, \bar{x}_0) = (\bar{x}, (A_0 - \lambda_0 E)\bar{x}_0) = 0,$$

так как $A_0 = A_0'$ и $(A_0 - \lambda_0 E)\bar{x}_0 = \bar{0}$. Таким образом, справедливость формулы (3.7.16) доказана. Аналогично проверяется и справедливость формулы (3.7.17), а именно:

$$\begin{aligned} (E - R(\mu E - B))\bar{x} &= \bar{x} - R(\mu E - B)\bar{x} = \bar{x} - R(A_0 - \lambda_0 E)\bar{x} = \\ &= \bar{x} - R(A_0 - \lambda_0 E)(\bar{x}_0 + \bar{z}) = \bar{x} - R(A_0 - \lambda_0 E)\bar{z} = \bar{x} - \bar{z} = \bar{x}_0, \end{aligned}$$

ибо по определению R при любом \bar{z} , ортогональном к \bar{x}_0 , $R(A_0 - \lambda_0 E)\bar{z} = \bar{z}$.

Пусть числа $\|B\|$ и $|\mu|$ настолько малы, что матрица $E - R(\mu E - B)$ невырожденная. Тогда из формулы (3.7.17) можно выразить вектор \bar{x} :

$$\bar{x} = (E - R(\mu E - B))^{-1} \bar{x}_0. \quad (3.7.18)$$

При малых $\|B\|$ и $|\mu|$ правую часть равенства (3.7.18) можно разложить в ряд. Ограничиваясь членами второго порядка малости относительно μ и B и принимая во внимание, что $R\bar{x}_0 = \bar{0}$, получим формулу для приближенного вычисления вектора \bar{x} :

$$\bar{x} \approx \bar{x}_0 - R B \bar{x}_0 + R B R B \bar{x}_0 - \mu R^2 B \bar{x}_0. \quad (3.7.19)$$

Умножив сначала слева это равенство на матрицу $\mu E - B$, а затем обе части скалярно на вектор \bar{x}_0 и принимая во внимание равенство (3.7.16), для приближенного вычисления μ найдем такую формулу:

$$\mu \approx \frac{(B\bar{x}_0, \bar{x}_0) - (R B \bar{x}_0, B \bar{x}_0) + (R B R B \bar{x}_0, R B \bar{x}_0)}{\|\bar{x}_0\|^2 + \|R B \bar{x}_0\|^2}. \quad (3.7.20)$$

Этим заканчивается решение задачи об определении вектора \bar{x} и числа μ .

Применим высказанные здесь соображения к решению задачи об уточнении отдельного собственного значения симметрической матрицы и отвечающего ему собственного вектора.

Предположим, что для рассматриваемой симметрической матрицы A мы каким-либо способом нашли приближенный собственный вектор \bar{x}_0 , который будем считать нормированным в евклидовой метрике, и приближенное собственное значение $\lambda_0 = (A\bar{x}_0, \bar{x}_0)$. Соответствующие точные значения для матрицы A обозначим через λ и \bar{x} . Полученные выше результаты можно будет применить, если мы построим такую симметрическую матрицу A_0 , для которой λ_0 и \bar{x}_0 будут точными собственными значением и вектором, т. е. будет выполняться условие $A_0\bar{x}_0 = \lambda_0\bar{x}_0$.

Можно показать, что искомой матрицей будет матрица

$$A_0 = A - \bar{x}_0\bar{r}' - \bar{r}\bar{x}_0',$$

где $\bar{r} = A\bar{x}_0 - \lambda_0\bar{x}_0$. Действительно, матрица A_0 симметрична; ибо

$$A_0' = A' - \bar{r}\bar{x}_0' - \bar{x}_0\bar{r}' = A - \bar{x}_0\bar{r}' - \bar{r}\bar{x}_0' = A_0$$

и для нее выполняется условие $A_0\bar{x}_0 = \lambda_0\bar{x}_0$, так как

$$\begin{aligned} A_0\bar{x}_0 &= \lambda_0\bar{x}_0 - (\bar{x}_0, \bar{r})\bar{x}_0 = \lambda_0\bar{x}_0 - (\bar{x}_0, A\bar{x}_0 - \lambda_0\bar{x}_0)\bar{x}_0 = \\ &= \lambda_0\bar{x}_0 - [(\bar{x}_0, A\bar{x}_0) - \lambda_0(\bar{x}_0, \bar{x}_0)]\bar{x}_0 = \lambda_0\bar{x}_0. \end{aligned}$$

Здесь выражение в квадратных скобках равно нулю, потому что по предположению $(\bar{x}_0, \bar{x}_0) = 1$ и $\lambda_0 = (A\bar{x}_0, \bar{x}_0)$.

На основании формул (3.7.19) и (3.7.20) теперь можно записать такие формулы для приближенного вычисления x и λ :

$$\bar{x} \approx \bar{x}_0 - R\bar{r} - \mu R^2\bar{r}, \quad (3.7.21)$$

$$\lambda \approx \lambda_0 - \frac{(R\bar{r}, \bar{r})}{1 + \|R\bar{r}\|^2}. \quad (3.7.22)$$

Здесь учтено, что

$$\|\bar{x}_0\|^2 = (\bar{x}_0, \bar{x}_0) = 1, \quad B\bar{x}_0 = (A - A_0)\bar{x}_0 = (\bar{x}_0\bar{r}' + \bar{r}\bar{x}_0')\bar{x}_0 = \bar{r}, \quad (\bar{r}, \bar{x}_0) = 0 \quad \text{и} \quad RB\bar{r} = 0.$$

При вычислениях по формулам (3.7.21) и (3.7.22) достаточно знать не матрицу R , а векторы $\bar{z}_1 = R\bar{r}$ и $\bar{z}_2 = R^2\bar{r}$, которые можно находить соответственно из систем

$$\left. \begin{aligned} (A_0 - \lambda_0 E)\bar{z}_1 &= \bar{r}, \\ (z_1, \bar{x}_0) &= 0 \end{aligned} \right\} \quad (3.7.23)$$

и

$$\left. \begin{aligned} (A_0 - \lambda_0 E)\bar{z}_2 &= \bar{z}_1, \\ (\bar{z}_2, \bar{x}_0) &= 0. \end{aligned} \right\} \quad (3.7.24)$$

Так как $|A_0 - \lambda_0 E| = 0$, то одно из первых n уравнений в системах (3.7.23) и (3.7.24) может быть отброшено при решении и затем использовано только лишь для контроля правильности вычислений. Найдя из этих систем векторы \bar{z}_1 и \bar{z}_2 , уточняем далее по формулам (3.7.21) и (3.7.22) собственный вектор и собственное значение.

Заметим, что если в формуле (3.7.21) отбросить член $\mu R^2 \bar{r}$ второго порядка малости относительно величин μ и \bar{r} , то отпадет необходимость в решении системы (3.7.24), из которой определяется вектор $\bar{z}_2 = R^2 \bar{r}$.

При проведении итерационного процесса посредством n -кратного повторного применения формул (3.7.21) и (3.7.22) скорость сходимости будет иметь порядок q^{5^n} , где $q = \frac{\|\bar{r}\|}{\tau - 2\|\bar{r}\|}$ и τ — расстояние от уточняемого собственного значения до ближай-

шего соседнего, $\bar{r} = A\bar{x}_0 - \lambda_0 \bar{x}_0$. Если аналогично использовать формулу (3.7.21) с отбрасыванием в ней члена $\mu R^2 \bar{r}$, то скорость сходимости будет иметь порядок q^{3^n} . Более подробно об оценках такого рода можно прочитать в упомянутой работе М. К. Гавурина [4].

3.7.3. δ^2 -Процесс Эйткена

В §3.5 мы отмечали, что степенной метод нахождения наибольшего по абсолютной величине собственного значения матрицы во многих случаях имеет недостаточно быструю сходимость, что требует выполнения большого числа итераций для получения решения с желаемой точностью. Ниже мы рассмотрим приемы ускорения сходимости некоторых последовательностей и, в частности, последовательностей, получающихся при использовании степенного метода. Одним из таких приемов является δ^2 -процесс Эйткена, о котором уже говорилось в § 1.3.

Напомним существо проблемы. Пусть задана числовая или функциональная последовательность $u_0, u_1, \dots, u_n, \dots$. Требуется преобразовать эту последовательность в новую последовательность $\{v_n\}$, которая сходилась бы к тому же самому пределу, что и последовательность $\{u_n\}$, но быстрее последней.

Каждый член последовательности $\{v_n\}$ будем определять по формуле Эйткена

$$v_n = \frac{u_{n+1}u_{n-1} - u_n^2}{u_{n+1} - 2u_n + u_{n-1}} \quad (n=1, 2, \dots), \quad (3.7.25)$$

где предполагается, что $u_{n+1} - 2u_n + u_{n-1} \neq 0$. Если последовательность $\{u_n\}$ сходится со скоростью геометрической прогрессии или близкой к ней, то, как указывалось в § 1.3, преобразование членов этой последовательности по формуле (3.7.25) может дать точное значение предела или сильно улучшить сходимость.

Если собственные значения матрицы A удовлетворяют условию

$$|\lambda_1| > |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|,$$

то преобразование Эйткена можно использовать для ускорения сходимости последовательности $\{\lambda_1^{(k)}\}$, возникающей в степенном методе. Известно, что в этом случае

$$\lambda_1^{(k)} = \lambda_1 + O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)$$

и

$$\lambda_1^{(k)} = \frac{y_s^{(k+1)}}{y_s^{(k)}}. \quad (3.7.26)$$

Поскольку $\lambda_1^{(k)}$ изменяется по закону, близкому к геометрической прогрессии, то можно ожидать, что последовательность $\{v_k\}$, каждый член которой определяется по формуле Эйткена

$$v_k = \frac{\lambda_1^{(k+1)} \lambda_1^{(k-1)} - [\lambda_1^{(k)}]^2}{\lambda_1^{(k+1)} - 2\lambda_1^{(k)} + \lambda_1^{(k-1)}} \quad (k=1, 2, \dots),$$

будет иметь более быструю сходимость к искомому пределу λ_1 — наибольшему по абсолютной величине собственному значению матрицы A .

Будем считать, что собственное значение λ_1 найдено достаточно точно. В этом случае δ^2 -процесс Эйткена можно применить также и к определению уточненного собственного вектора, отвечающего λ_1 . Приближенно собственный вектор, отвечающий λ_1 , в силу формулы (3.5.8) равен $\bar{y}^{(k)} = A^{(k)} \bar{y}^{(0)}$. Правило уточнения собственного вектора построим таким образом, чтобы каждая компонента вектора $\bar{y}^{(k)}$ уточнялась отдельно. В силу формулы (3.5.4) для s -й компоненты вектора $\bar{y}^{(k)}$ имеем

$$y_s^{(k)} = \beta_{s1} \lambda_1^k + \beta_{s2} \lambda_2^k + \dots + \beta_{sn} \lambda_n^k.$$

Рассмотрим наряду с вектором $\bar{y}^{(k)}$ векторы $\bar{y}^{(k-1)}$, $\bar{y}^{(k+1)}$ и выберем в этих векторах компоненты $y_s^{(k-1)}$, $y_s^{(k+1)}$. Составим величины

$$\lambda_1 \cdot y_s^{(k-1)}, \quad 1 \cdot y_s^{(k)}, \quad -\frac{1}{\lambda_1} \cdot y_s^{(k+1)}$$

и применим к ним формулу Эйткена. Тогда получим

$$\begin{aligned} v_s^{(k)} &= \frac{-\frac{1}{\lambda_1} y_s^{(k+1)} \lambda_1 y_s^{(k-1)} - [y_s^{(k)}]^2}{-\frac{1}{\lambda_1} y_s^{(k+1)} - 2y_s^{(k)} + \lambda_1 y_s^{(k-1)}} = \\ &= \beta_{s1} \lambda_1^k \left[1 + O\left(\left|\frac{\lambda_3}{\lambda_1}\right|^k\right) \right]. \end{aligned} \quad (3.7.27)$$

Аналогичное выражение можно записать и для $y_s^{(k)}$:

$$y_s^{(k)} = \beta_{s1} \lambda_1^k \left[1 + O \left(\left| \frac{\lambda_2}{\lambda_1} \right|^k \right) \right]. \quad (3.7.28)$$

Поскольку по предположению $\left| \frac{\lambda_3}{\lambda_1} \right| < \left| \frac{\lambda_2}{\lambda_1} \right|$, то из формул (3.7.27) и (3.7.28) следует, что последовательность $\{v_s^{(k)}\}$ будет быстрее, чем последовательность $\{y_s^{(k)}\}$, сходиться к искомому пределу — s -й компоненте собственного вектора, отвечающего λ_1 . Причем эта сходимость будет тем быстрее, чем $|\lambda_3|$ меньше $|\lambda_2|$.

Укажем еще, как в некоторых случаях изложенный здесь прием можно использовать и для решения задачи об ускорении сходимости итерационных последовательностей, возникающих при решении систем уравнений. Пусть рассматривается следующая система:

$$A\bar{x} = \bar{b}.$$

Запишем для нее формулу стационарного линейного итерационного процесса в виде

$$\bar{x}^{(k+1)} = B\bar{x}^{(k)} + C\bar{b}, \quad (3.7.29)$$

где B и C — такие матрицы, что

$$B + CA = E. \quad (3.7.30)$$

Будем считать, что матрица B имеет линейные элементарные делители и все ее собственные значения вещественны, причем $|\lambda_i| < 1$ ($i = 1, 2, \dots, n$). Предположим также, что λ_1 является наибольшим по абсолютной величине собственным значением матрицы B . Обозначим собственные векторы этой матрицы через $\bar{z}_1, \bar{z}_2, \dots, \bar{z}_n$ и разложим по ним вектор $\bar{x}^{(*)} - \bar{x}^{(0)}$, где $\bar{x}^{(*)}$ — точное решение системы $A\bar{x} = \bar{b}$, а $\bar{x}^{(0)}$ — некоторое начальное приближение к решению:

$$\bar{x}^{(*)} - \bar{x}^{(0)} = \alpha_1 \bar{z}_1 + \alpha_2 \bar{z}_2 + \dots + \alpha_n \bar{z}_n. \quad (3.7.31)$$

Как ранее было показано (п. 2.2.2), имеет место равенство

$$\bar{x}^{(*)} - \bar{x}^{(m)} = B^m (\bar{x}^{(*)} - \bar{x}^{(0)}), \quad (3.7.32)$$

где $\bar{x}^{(m)}$ — приближение номера m , полученное по формуле (3.7.29) при начальном векторе $\bar{x}^{(0)}$. В силу формул (3.7.31) и (3.7.32) имеем

$$\bar{x}^{(*)} - \bar{x}^{(m)} = \alpha_1 \lambda_1^m \bar{z}_1 + \alpha_2 \lambda_2^m \bar{z}_2 + \dots + \alpha_n \lambda_n^m \bar{z}_n. \quad (3.7.33)$$

Записывая это равенство для s -й компоненты векторов, получим

$$x_s^{(*)} - x_s^{(m)} = \alpha_1 \lambda_1^m z_{s1} + \alpha_2 \lambda_2^m z_{s2} + \dots + \alpha_n \lambda_n^m z_{sn}, \quad (3.7.34)$$

где z_{sk} есть s -я компонента вектора \bar{z}_k , $s = 1, 2, \dots, n$. Формулу (3.7.34) можно записать иначе:

$$x_s^{(m)} = x_s^{(*)} - \alpha_1 \lambda_1^m z_{s1} \left[1 + O \left(\left| \frac{\lambda_2}{\lambda_1} \right|^m \right) \right],$$

в предположении, что m достаточно велико. Из этой формулы видно, что $x_s^{(m)}$ изменяется по закону, близкому к геометрической прогрессии, поэтому последовательность $\{v_s^{(m)}\}$, где

$$v_s^{(m)} = \frac{x_s^{(m+1)} x_s^{(m-1)} - (x_s^{(m)})^2}{x_s^{(m+1)} - 2x_s^{(m)} + x_s^{(m-1)}},$$

будет быстрее сходиться к $x_s^{(*)}$ — s -й компоненте искомого решения $\bar{x}^{(*)}$, чем последовательность $\{x_s^{(m)}\}$.

§2. Процесс Эйткена можно обобщить и на итерационные процессы с матрицей B , у которой преобладающих по абсолютной величине собственных значений будет два или больше. В этом случае формулы Эйткена усложняются, и мы их здесь приводить не будем. Они имеются в учебном пособии [1].

3.7.4. Метод М. К. Гавурина

Пусть, как и ранее, рассматривается система линейных алгебраических уравнений $A\bar{x} = \bar{b}$ и для численного нахождения решения этой системы избрана формула (3.7.29) с некоторым начальным вектором $\bar{x}^{(0)}$. Предположим, что матрица B удовлетворяет сформулированным в предыдущем пункте условиям и итерационная последовательность, построенная по формуле (3.7.29), имеет медленную сходимость к своему пределу $\bar{x}^{(*)}$ — решению системы $A\bar{x} = \bar{b}$. Это будет иметь место, например, в случае, когда все собственные значения λ_i матрицы B по абсолютной величине меньше единицы, но среди $|\lambda_i|$ есть некоторые, близкие к единице. Действительно, в силу формулы (3.7.33) при некотором $|\lambda_i| \approx 1$ в векторе ошибок $\bar{x}^{(*)} - \bar{x}^{(m)}$ составляющая $\alpha_i \lambda_i^m \bar{x}_i$ будет медленно убывать при воз-

растании m . Как можно ускорить сходимость в таком случае? Очевидно, что решению проблемы способствовали бы такие преобразования последовательности $\{\bar{x}^{(m)}\}$, которые позволили бы уменьшить влияние компоненты $\alpha_i \lambda_i^m \bar{x}_i$ и некоторых других компонент, для которых $|\lambda_j|$ близок к единице, в разложениях типа (3.7.33).

Делается это следующим путем. Пусть, начиная с некоторого вектора $\bar{x}^{(0)}$, мы вычислили по формуле (3.7.29) векторы $\bar{x}^{(1)}, \bar{x}^{(2)}, \dots, \bar{x}^{(s+1)}$. Составим разности $\bar{x}^{(1)} - \bar{x}^{(0)}, \bar{x}^{(2)} - \bar{x}^{(1)}, \dots, \bar{x}^{(s+1)} - \bar{x}^{(s)}$ и укажем правило для определения некоторых коэффициентов α_i таким образом, чтобы вектор \bar{y} , вычисляемый по формуле

$$\bar{y} = \bar{x}^{(0)} + \alpha_0(\bar{x}^{(1)} - \bar{x}^{(0)}) + \dots + \alpha_s(\bar{x}^{(s+1)} - \bar{x}^{(s)}), \quad (3.7.35)$$

возможно точнее совпал с вектором $\bar{x}^{(*)}$ — решением системы $A\bar{x} = \bar{b}$. Запишем разложение вектора $\bar{x}^{(1)} - \bar{x}^{(0)}$ по собственным векторам \bar{z}_i матрицы B :

$$\bar{x}^{(1)} - \bar{x}^{(0)} = \beta_1 \bar{z}_1 + \beta_2 \bar{z}_2 + \dots + \beta_n \bar{z}_n.$$

Используя формулу (3.7.29), для вектора $\bar{x}^{(2)} - \bar{x}^{(1)}$ получим

$$\bar{x}^{(2)} - \bar{x}^{(1)} = B(\bar{x}^{(1)} - \bar{x}^{(0)}).$$

Значит,

$$\bar{x}^{(2)} - \bar{x}^{(1)} = \beta_1 \lambda_1 \bar{z}_1 + \beta_2 \lambda_2 \bar{z}_2 + \dots + \beta_n \lambda_n \bar{z}_n.$$

По аналогии для любого вектора $\bar{x}^{(k+1)} - \bar{x}^{(k)} = B(\bar{x}^{(k)} - \bar{x}^{(k-1)})$ найдем выражение

$$\begin{aligned} \bar{x}^{(k+1)} - \bar{x}^{(k)} &= \beta_1 \lambda_1^k \bar{z}_1 + \beta_2 \lambda_2^k \bar{z}_2 + \dots + \beta_n \lambda_n^k \bar{z}_n = \sum_{i=1}^n \beta_i \lambda_i^k \bar{z}_i \quad (3.7.36) \\ (k &= 0, 1, 2, \dots). \end{aligned}$$

По условию последовательность $\{\bar{x}^{(k)}\}$ — сходящаяся ($\lim_{k \rightarrow \infty} \bar{x}^{(k)} = \bar{x}^{(*)}$), поэтому точное решение $\bar{x}^{(*)}$ можно записать в виде

$$\begin{aligned} \bar{x}^{(*)} &= \bar{x}^{(0)} + \sum_{k=0}^{\infty} (\bar{x}^{(k+1)} - \bar{x}^{(k)}) = \bar{x}^{(0)} + \sum_{k=0}^{\infty} \sum_{i=1}^n \beta_i \lambda_i^k \bar{z}_i = \\ &= \bar{x}^{(0)} + \sum_{i=1}^n \frac{\beta_i}{1 - \lambda_i} \bar{z}_i, \end{aligned} \quad (3.7.37)$$

ибо

$$|\lambda_i| < 1 \quad \text{и} \quad \sum_{k=0}^{\infty} \lambda_i^k = \frac{1}{1-\lambda_i}.$$

Для вектора же \bar{y} в силу формул (3.7.35) и (3.7.36) получим теперь такое выражение:

$$\bar{y} = \bar{x}^{(0)} + \sum_{k=0}^s \alpha_k \sum_{i=1}^n \beta_i \lambda_i^k \bar{z}_i = \bar{x}^{(0)} + \sum_{i=1}^n \beta_i P(\lambda_i) \bar{z}_i, \quad (3.7.38)$$

где через $P(\lambda)$ обозначен многочлен

$$P(\lambda) = \sum_{k=0}^s \alpha_k \lambda^k.$$

Таким образом, разность $\bar{x}^{(*)} - \bar{y}$ можно представить в форме

$$\bar{x}^{(*)} - \bar{y} = \sum_{i=1}^n \beta_i \left[\frac{1}{1-\lambda_i} - P(\lambda_i) \right] \bar{z}_i. \quad (3.7.39)$$

По смыслу задачи мы должны стремиться уменьшить разность $\bar{x}^{(*)} - \bar{y}$. Сделать это мы сможем за счет выбора коэффициентов α_i , ибо из формулы (3.7.39) видно, что чем меньше будут разности

$$\frac{1}{1-\lambda_i} - P(\lambda_i)$$

по модулю, тем меньше будет и норма вектора $\bar{x}^{(*)} - \bar{y}$.

Пусть $\max_{1 \leq i \leq n} |\lambda_i| = M < 1$. В связи с вышеизложенным, задача минимизации разности $\bar{x}^{(*)} - \bar{y}$ сводится теперь к следующему: найти многочлен $P(\lambda)$ степени не выше s , такой, чтобы он на отрезке $[-M, M]$ наименее отклонялся от функции $\frac{1}{1-\lambda}$. Полином степени s , удовлетворяющий последнему требованию, известен (см. [1], гл. IV, задача 10) и равен

$$P(\lambda) = \frac{2\alpha^{s+1}}{(1-\alpha^2)^2} \cdot \frac{1}{\lambda-1} \left[T_{s+1} \left(\frac{\lambda}{M} \right) - 2\alpha T_s \left(\frac{\lambda}{M} \right) + \right. \\ \left. + \alpha^2 T_{s-1} \left(\frac{\lambda}{M} \right) \right] + \frac{1}{\lambda-1},$$

где

$$\alpha = \frac{1}{M} - \sqrt{\frac{1}{M^2} - 1}, \quad T_{i+1}(t) = 2tT_i(t) - T_{i-1}(t),$$

$$T_0(t) = 1, \quad T_1(t) = t, \quad i = 1, 2, \dots, s, \quad t \in [-1, 1].$$

Следовательно, алгоритм метода Гавурина может быть реализован таким образом:

- 1) по формуле (3.7.29) вычисляем несколько приближений к решению, например, $\bar{x}^{(1)}, \bar{x}^{(2)}, \dots, \bar{x}^{(s+1)}$ ($s \geq 1$);
- 2) определяем каким-либо приближенным путем число M и по указанной выше формуле находим коэффициенты α_k многочлена $P(\lambda)$;
- 3) используя эти коэффициенты, получаем по формуле (3.7.35) вектор \bar{y} , являющийся более точным приближением к искомому решению системы, чем вектор $\bar{x}^{(s+1)}$.

Отметим, что для больших значений s вычисления по формуле (3.7.35) могут быть связаны с потерей значащих цифр из-за сильного возрастания коэффициентов α_k . Поэтому s большим брать в разложении (3.7.35) не всегда можно.

3.7.5. Метод Л. А. Люстерника

Идея метода заключается в выделении главной части из остатка. Пусть опять мы рассматриваем систему уравнений $A\bar{x} = \bar{b}$ и решение этой системы отыскиваем по итерационной формуле (3.7.29). Предположим, что мы вычислили несколько членов итерационной последовательности $\{\bar{x}^{(k)}\}$ и убедились, что ее сходимость недостаточно быстрая. Поэтому перед нами встает вопрос об ускорении сходимости последовательности $\{\bar{x}^{(k)}\}$. Это можно сделать в отдельных случаях, а именно, когда:

- 1) матрица B обладает полной системой собственных векторов;
- 2) среди собственных чисел λ_i матрицы B есть наибольшее по модулю, т. е.

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|. \quad (3.7.40)$$

Поскольку мы предположили последовательность $\{\bar{x}^{(k)}\}$ сходящейся к решению $\bar{x}^{(*)}$ системы $A\bar{x} = \bar{b}$ и начальный вектор $\bar{x}^{(0)}$ считается произвольным, то мы должны считать, что $|\lambda_i| < 1$ ($i = 1, 2, \dots, n$).

Итак, получим остаток $\bar{x}^{(*)} - \bar{x}^{(k)}$ и найдем его главную часть. По формуле (3.7.36) имеем

$$\bar{x}^{(k+1)} - \bar{x}^{(k)} = \sum_{i=1}^n \beta_i \lambda_i^k \bar{z}_i \quad (k = 0, 1, 2, \dots). \quad (3.7.41)$$

Значит,

$$\bar{x}^{(*)} - \bar{x}^{(k)} = \sum_{i=0}^{\infty} (\bar{x}^{(i+k+1)} - \bar{x}^{(i+k)}) = \sum_{j=1}^n \frac{\lambda_j^k}{1 - \lambda_j} \beta_j \bar{z}_j. \quad (3.7.42)$$

Если k достаточно велико, то в силу условий (3.7.40) из формул (3.7.41) и (3.7.42) сможем выделить главные части и получить приближенные равенства

$$\bar{x}^{(k+1)} - \bar{x}^{(k)} \approx \beta_1 \lambda_1^k \bar{z}_1, \quad (3.7.43)$$

$$\bar{x}^{(*)} - \bar{x}^{(k)} \approx \beta_1 \frac{\lambda_1^k}{1 - \lambda_1} \bar{z}_1. \quad (3.7.44)$$

Определим из формулы (3.7.43) вектор $\beta_1 \bar{z}_1$ и подставим это значение в формулу (3.7.44), тогда для вектора $\bar{x}^{(*)}$ приближенно получим

$$\bar{x}^{(*)} \approx \bar{x}^{(k)} + \frac{1}{1 - \lambda_1} (\bar{x}^{(k+1)} - \bar{x}^{(k)}).$$

Следовательно, можно ожидать, что вектор \bar{y} , определяемый по формуле

$$\bar{y} = \bar{x}^{(k)} + \frac{1}{1 - \tilde{\lambda}_1} (\bar{x}^{(k+1)} - \bar{x}^{(k)}), \quad (3.7.45)$$

где через $\tilde{\lambda}_1$ обозначено приближенное наибольшее по модулю собственное значение матрицы B , будет ближе к $\bar{x}^{(*)}$, чем вектор $\bar{x}^{(k)}$ или $\bar{x}^{(k+1)}$.

Оценим разность $\bar{x}^{(*)} - \bar{y}$. Будем считать, что для вычисления $\tilde{\lambda}_1$ используется последовательность $\{\bar{x}^{(k)}\}$. Это возможно осуществить, ибо в силу формулы (3.7.41) имеем:

$$\begin{aligned} \bar{x}^{(k+1)} - \bar{x}^{(k)} &= \beta_1 \lambda_1^k \bar{z}_1 + \beta_2 \lambda_2^k \bar{z}_2 + \dots + \beta_n \lambda_n^k \bar{z}_n, \\ \bar{x}^{(k)} - \bar{x}^{(k-1)} &= \beta_1 \lambda_1^{k-1} \bar{z}_1 + \beta_2 \lambda_2^{k-1} \bar{z}_2 + \dots + \beta_n \lambda_n^{k-1} \bar{z}_n. \end{aligned}$$

Значит, $\tilde{\lambda}_1$ может быть определено как отношение одноименных компонент векторов $\bar{x}^{(k+1)} - \bar{x}^{(k)}$ и $\bar{x}^{(k)} - \bar{x}^{(k-1)}$. В этом случае, как это было показано в п. 3.5.1, можно положить

$$\lambda_1 = \tilde{\lambda}_1 + \varepsilon,$$

где

$$\varepsilon = O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right).$$

Установим связь между векторами $\bar{x}^{(*)} - \bar{y}$ и $\bar{x}^{(*)} - \bar{x}^{(k)}$, для чего введем в рассмотрение матрицу B_1 , определяемую равенством

$$B_1 = \frac{1}{1 - \tilde{\lambda}_1} [B - \tilde{\lambda}_1 E].$$

Легко проверить, что $\bar{x}^{(*)} - \bar{y} = B_1(\bar{x}^{(*)} - \bar{x}^{(k)})$. Действительно,

$$\begin{aligned} \bar{x}^{(*)} - \bar{y} &= \bar{x}^{(*)} - \bar{x}^{(k)} - \frac{1}{1 - \tilde{\lambda}_1} (\bar{x}^{(k+1)} - \bar{x}^{(k)}) = \bar{x}^{(*)} - \bar{x}^{(k)} - \\ &- \frac{1}{1 - \tilde{\lambda}_1} [(\bar{x}^{(k+1)} - \bar{x}^{(*)}) + (\bar{x}^{(*)} - \bar{x}^{(k)})] = \frac{1}{1 - \tilde{\lambda}_1} B(\bar{x}^{(*)} - \bar{x}^{(k)}) + \\ &+ \left(1 - \frac{1}{1 - \tilde{\lambda}_1}\right) (\bar{x}^{(*)} - \bar{x}^{(k)}) = B_1(\bar{x}^{(*)} - \bar{x}^{(k)}), \end{aligned}$$

так как $\bar{x}^{(*)} - \bar{x}^{(k+1)} = B(\bar{x}^{(*)} - \bar{x}^{(k)})$. Далее, используя формулу (3.7.42), получим

$$\begin{aligned} \bar{x}^{(*)} - \bar{y} &= B_1(\bar{x}^{(*)} - \bar{x}^{(k)}) = \frac{1}{1 - \tilde{\lambda}_1} \left[\frac{\lambda_1^k \varepsilon}{1 - \lambda_1} \beta_1 \bar{z}_1 + \right. \\ &+ \frac{\lambda_2^k (\lambda_2 - \tilde{\lambda}_1)}{1 - \lambda_2} \beta_2 \bar{z}_2 + \dots + \frac{\lambda_n^k (\lambda_n - \tilde{\lambda}_1)}{1 - \lambda_n} \beta_n \bar{z}_n \left. \right]. \end{aligned}$$

Учитывая, что $\varepsilon = O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)$, будем иметь

$$\bar{x}^{(*)} - \bar{y} = O(|\lambda_2|^k) \bar{z}_1. \quad (3.7.46)$$

Так как в силу формулы (3.7.44)

$$\bar{x}^{(*)} - \bar{x}^{(k)} = O(|\lambda_1|^k) \bar{z}_1, \quad (3.7.47)$$

то из равенств (3.7.46) и (3.7.47) следует, что улучшение сходимости будет тем больше, чем меньше отношение $\left|\frac{\lambda_2}{\lambda_1}\right|$.

Если $\tilde{\lambda}_1$ близко к 1, то множитель $\frac{1}{1-\tilde{\lambda}_1}$ может принимать большое значение и при вычислениях по формуле (3.7.45) из-за этого может произойти потеря значащих цифр. Поэтому вместо названной формулы целесообразно пользоваться следующей:

$$\bar{y} = \bar{x}^{(k)} + \frac{1}{1-\tilde{\lambda}_1^p} (\bar{x}^{(k+p)} - \bar{x}^{(p)})$$

(при условии, что $\tilde{\lambda}_1^p$ существенно меньше 1). Эта формула выводится так же, как и формула (3.7.45).

Метод Л. А. Люстерника можно распространить и на случай, когда матрица B имеет несколько равных наибольших по модулю собственных значений, т. е.

$$|\lambda_1| = |\lambda_2| = \dots = |\lambda_r| > |\lambda_{r+1}| \geq \dots \geq |\lambda_n|, \quad r \geq 2.$$

Покажем, как это можно сделать при наличии у матрицы пары комплексно сопряженных собственных значений, наибольших по модулю. Матрицу B будем считать вещественной и предположим, что

$$|\lambda_1| = |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|$$

и

$$\lambda_1 = \bar{\lambda}_2,$$

где через $\bar{\lambda}_2$ обозначено число, комплексно сопряженное с λ_2 . Собственным значениям λ_1 и λ_2 соответствуют комплексно сопряженные собственные векторы \bar{z}_1 и \bar{z}_2 матрицы B . При достаточно большом k из формул (3.7.41) и (3.7.42) получим приближенные равенства

$$\bar{x}^{(k+1)} - \bar{x}^{(k)} \approx \beta_1 \lambda_1^k \bar{z}_1 + \beta_2 \lambda_2^k \bar{z}_2, \quad (3.7.48)$$

$$\bar{x}^{(k+2)} - \bar{x}^{(k+1)} \approx \beta_1 \lambda_1^{k+1} \bar{z}_1 + \beta_2 \lambda_2^{k+1} \bar{z}_2 \quad (3.7.49)$$

и

$$\bar{x}^{(*)} - \bar{x}^{(k)} \approx \frac{\lambda_1^k}{1-\lambda_1} \beta_1 \bar{z}_1 + \frac{\lambda_2^k}{1-\lambda_2} \beta_2 \bar{z}_2. \quad (3.7.50)$$

Из равенств (3.7.48) и (3.7.49) найдем векторы $\beta_1 \bar{z}_1$ и $\beta_2 \bar{z}_2$ и подставим эти значения в равенство (3.7.50). Тогда для вектора $\bar{x}^{(*)}$ будет справедливо следующее приближенное равенство:

$$\bar{x}^{(*)} \approx \bar{x}^{(k)} + \frac{[1 - (\lambda_1 + \lambda_2)] \bar{p} + \bar{s}}{1 - (\lambda_1 + \lambda_2) + \lambda_1 \lambda_2},$$

где $\bar{s} = \bar{x}^{(k+2)} - \bar{x}^{(k+1)}$, $\bar{p} = \bar{x}^{(k+1)} - \bar{x}^{(k)}$.

Таким образом, в качестве вектора \bar{y} , улучшенного приближения к $\bar{x}^{(*)}$, в этом случае следует взять вектор

$$\bar{y} = \bar{x}^{(k)} + \frac{[1 - (\tilde{\lambda}_1 + \tilde{\lambda}_2)] \bar{p} + \bar{s}}{1 - (\tilde{\lambda}_1 + \tilde{\lambda}_2) + \tilde{\lambda}_1 \tilde{\lambda}_2},$$

где через $\tilde{\lambda}_1$ и $\tilde{\lambda}_2$ обозначена пара приближенных комплексно сопряженных собственных значений матрицы B .

Литература

1. Березин И. С., Жидков Н. П. Методы вычислений, т. I. М., 1966.
2. Воеводин В. В. Численные методы алгебры (теория и алгоритмы). М., 1966.
3. Гавурин М. К. Применение полиномов наилучшего приближения к улучшению сходимости итерационных процессов. УМН, 5: 3, (37), 1956.
4. Гавурин М. К. О методе ложных возмущений для разыскания собственных значений. Журн. вычисл. матем. и матем. физики, I, № 5, 1961.
5. Демидович Б. П., Марон И. А. Основы вычислительной математики. М., 1963.
6. Ланцош К. Практические методы прикладного анализа. М., 1961.
7. Люстерник Л. А. Замечания к численному решению краевых задач для уравнения Лапласа и вычислению собственных значений методом сеток. Тр. матем. ин-та им. В. А. Стеклова, 20, 1947.
8. Милн В. Э. Численный анализ. М., 1951.
9. Фаддеев Д. К., Фаддеева В. Н. Вычислительные методы линейной алгебры. М., 1963.
10. Хаусхолдер А. С. Основы численного анализа. М., 1956.
11. Aitken A. Studies in practical mathematics, II. The evaluation of the latent roots and latent vectors of a matrix. Proc. Roy. Soc. Edinburgh, Sec. A., 1936, 1937, 57.
12. Derwidue L. Une methode mecanique de calcul des vecteurs d'une matrice quelconque. Bull. Soc. roy. sci. hiege, 24, № 5, 1955.

Глава 4

ИНТЕРПОЛИРОВАНИЕ

§ 4.1. О СОДЕРЖАНИИ ЗАДАЧИ ИНТЕРПОЛИРОВАНИЯ

4.1.1. Об интерполяционных приближениях

Слово «интерполирование» означает нахождение промежуточных значений. В математике этому термину придают более определенный, однако не всегда одинаковый смысл. Причиной такого расхождения в понимании в первую очередь является то обстоятельство, что проблемы, которые можно отнести к интерполяционным, очень разнообразны и методы решения их могут сильно различаться между собой.

Мы будем рассматривать лишь частную задачу интерполирования, в которой целью является нахождение значений функции. В достаточно общей форме эта задача может быть высказана в следующих словах. Пусть на конечном или бесконечном отрезке $\langle a, b \rangle$ рассматривается некоторая m -кратно непрерывно дифференцируемая функция f . Предположим, что в k_0 точках $x_{01}, x_{02}, \dots, x_{0k_0}$ известны значения функции $f(x_{01}), f(x_{02}), \dots, f(x_{0k_0})$, в k_1 точках $x_{11}, x_{12}, \dots, x_{1k_1}$ известны значения первой производной от нее $f'(x_{11}), \dots, f'(x_{1k_1})$ и т. д. и в k_m точках x_{m1}, \dots, x_{mk_m} известны значения производной от f порядка m : $f^{(m)}(x_{m1}), \dots, f^{(m)}(x_{mk_m})$.

Все перечисленные величины называются исходными данными интерполирования, а точки x_{ij} — узлами интерполирования.

Общее число известных значений функции f и производных обозначим n : $k_0 + k_1 + \dots + k_m = n$.

Возьмем любую точку $x \in \langle a, b \rangle$, отличную от узлов x_{0i} ($i = 1, \dots, k_0$), и поставим себе задачу найти, пользуясь исходными данными, значение $f(x)$. Такая задача является, очевидно, весьма неопределенной и может решаться лишь приближенно. Чтобы понять степень ее неопределенности, достаточно рассмотреть простейший случай задачи, когда интерполирование выполняется по значениям только самой функции, значения же производных отсутствуют. Пусть в n точках x_1, x_2, \dots, x_n известны значения $y_1 = f(x_1), y_2 = f(x_2), \dots, y_n = f(x_n)$ функции f . Геометрически это означает, что в плоскости с системой координат xu даны n точек $M_k(x_k, y_k)$ ($k = 1, 2, \dots, n$) на графике l функции f . По ним мы должны найти $f(x)$, и так как x может быть любой точкой отрезка $\langle a, b \rangle$, то дело идет о восстановлении графика $f(x)$ на $\langle a, b \rangle$.

Прежде всего необходимо условиться о том, должна ли линия, которую мы предполагаем приближенно принять за график f , точно прохо-

дить через заданные точки M_k или она должна идти лишь достаточно близко от них.

Необходимость в таком условии возникает, например, в том случае, когда значения $f(x_k)$ находятся из опыта. Точность таких значений бывает, как правило, невысокой и ограничивается точностью измерений. В таких задачах излишне требовать, чтобы линия, которую мы должны построить, проходила точно через точки M_k , достаточно, чтобы она отклонялась от M_k по ординате y на величину, не большую погрешности измерений.

Такой же вопрос, если говорить о принципиальной стороне дела, возникает и при нахождении значений функции f , заданной таблично. Точность совпадения $f(x_k)$ и ординат нужной нам линии при $x=x_k$ ($k=1, \dots, n$) должна быть не выше точности таблицы.

В нашем изложении мы встанем на крайнюю точку зрения и будем требовать точного прохождения линии через M_k . Отметим лишь, что такая постановка задачи интерполирования является не единственно возможной и не во всех случаях самой целесообразной.

Допустим, что через точки M_k ($k=1, \dots, n$) мы провели некоторую линию λ и пусть $\varphi(x)$ есть функция, для которой λ будет графиком. $\varphi(x)$ ниже мы будем называть интерполяционным приближением к $f(x)$. По построению φ будет выполнять равенства

$$\varphi(x_k) = f(x_k) \quad (k=1, 2, \dots, n). \quad (4.1.1)$$

Таких линий λ и соответствующих им функций φ существует бесконечное множество.

В узлах интерполирования f и φ совпадают, но, когда x изменяется между узлами или в стороне от их расположения, линии f и λ могут расходиться и расхождение может быть весьма сильным, даже если узлов много и они мало удалены друг от друга. Так может случиться, например, когда f будет сильно извилистой или даже разрывной линией, а в качестве λ мы возьмем аналитическую линию или линию, обладающую высоким порядком гладкости. Чтобы надеяться получить удовлетворительное совпадение f и λ , нужно как-то согласовать между собой свойства этих линий. Но выбор способа проведения линии λ есть не что иное, как выбор правила интерполирования. Если с наглядного геометрического языка перейти на язык числовых переменных, то можно сказать, что для получения удовлетворительного правила интерполирования необходимо способ интерполирования в каких-то разумных границах согласовать с заранее известными свойствами интерполируемой функции f , такими, как непрерывность, дифференцируемость, аналитичность и др. О таком согласовании мы будем говорить более подробно немного позже, а сейчас приведем простые пояснительные примеры. Если функция f достаточно гладкая и ее нужно интерполировать на конечном отрезке, то можно надеяться получить хорошую точность, если интерполирующую функцию φ искать среди целых алгебраических многочленов.

Когда необходимо интерполировать функцию f на полуоси $[0, \infty)$, гладкую там и имеющую конечный предел $\lim_{x \rightarrow \infty} f(x) = f(+\infty)$, то в качестве интерполирующей функции можно взять рациональную функцию

$$\varphi(x) = \frac{P(x)}{Q(x)},$$

где P и Q — многочлены, причем степень P не больше степени Q и нули знаменателя Q лежат вне полуоси $[0, \infty)$.

Наконец, если мы интерполируем на всей оси гладкую периодическую функцию, то можем воспользоваться тригонометрическим многочленом с тем же периодом и т. д.

В предшествующем изложении мы хотели обратить внимание на то, что в задаче интерполирования функций, даже при строгом выполнении исходных условий типа (4.1.1), остается еще большой произвол в способе построения интерполяционного приближения φ и этим произволом нужно воспользоваться так, чтобы добиться, насколько это окажется возможным, лучшей точности результата в нахождении значений f .

Возвратимся теперь к общей задаче, сформулированной нами раньше, об интерполировании $f(x)$ по нескольким значениям самой функции и производных от неё до порядка m .

Пусть задано множество F функций f , подлежащих интерполированию. Обычно это бывает множество функций, обладающих одинаковыми структурными свойствами, такими, например, как одинаковый порядок дифференцируемости, периодичность, аналитичность и др. Предположим также, что мы выбрали множество Φ функций φ , среди которых будем находить интерполяционные приближения для каждой функции f . Функция φ должна удовлетворять условиям

$$\varphi^{(j)}(x_{ij}) = f^{(j)}(x_{ij}) \quad (j=1, 2, \dots, k_i, i=0, 1, \dots, m). \quad (4.1.2)$$

Нашей ближайшей целью будет указать на некоторые требования, которым должен быть подчинен выбор множества Φ , чтобы при интерполировании можно было надеяться на хорошую точность.

Прежде всего необходимо указать на техническое требование, которому должны удовлетворять функции φ . Они применяются для приближения и нахождения значений f и поэтому должны быть достаточно простыми и удобными для вычислений.

Число m , являющееся наивысшим порядком производных, входящих в условия (4.1.2), будем считать для простоты фиксированным, что же касается чисел k_i ($i=0, 1, \dots, m$) и, следовательно, общего числа условий $n=k_0+k_1+\dots+k_m$, то их мы будем предполагать произвольными, но определенными в каждой частной задаче интерполирования.

Формальные требования, которые предъявляет к множеству Φ сама проблема интерполирования, состоят в следующем.

Во-первых, в условия (4.1.2) входят значения $\varphi^{(m)}(x_{mj})$ и, так как узлы x_{mj} ($j=1, \dots, k_m$) могут лежать в любом месте отрезка $\langle a, b \rangle$ и f непрерывно дифференцируема там, функции φ должны быть m -кратно непрерывно дифференцируемыми на $\langle a, b \rangle$.

Во-вторых, число условий (4.1.2) равно n и, чтобы выполнить их, нужно, вообще говоря, иметь семейство функций φ , зависящее от n численных параметров:

$$\varphi = \varphi_n(x, a_1, a_2, \dots, a_n). \quad (4.1.3)$$

Последние выбирают так, чтобы выполнялись уравнения

$$\frac{\partial^i}{\partial x^i} \varphi_n(x_{ij}, a_1, a_2, \dots, a_n) = f^{(i)}(x_{ij}). \quad (4.1.4)$$

$$(j=1, 2, \dots, k_i; i=0, 1, \dots, m).$$

В реальных задачах множества функций f , для которых строятся правила интерполирования, бывают весьма широкими и среди функций f будет существовать такая, для которой $f^{(i)}(x_{ij})$ будут равны наперед заданным произвольным числам. Поэтому необходимо требовать, чтобы система (4.1.4) была разрешимой при любых правых частях $f^{(i)}(x_{ij})$.*)

Третье требование, предъявляемое к выбору множества Φ , которое будет определено через несколько строк, является, по сути дела, условием плотности Φ повсюду в множестве F . Оно не достаточно для возможности сколь угодно точного интерполирования f , но, как будет видно из дальнейшего изложения, является необходимым в том смысле, что если это условие не выполняется, то интерполировать функцию $f \in F$ со сколь угодно малой погрешностью можно будет лишь в исключительных, как правило, случаях.

До настоящего места мы предполагали узлы x_{ij} фиксированными и говорили только об одном шаге интерполяционного процесса. Он определяется таблицей узлов **)

*) По самому смыслу интерполяционной проблемы, если заданы исходные условия (4.1.2), то интерполяционное приближение φ должно определяться однозначно. Когда система (4.1.4) имеет несколько решений, то при подстановке их в (4.1.3) должна получиться одна и та же функция φ_n . Если же разным решениям системы будут отвечать разные функции φ_n , то выбор семейства (4.1.3) следует, по-видимому, признать неудачным и отказаться строить при помощи такой $\varphi_n(x, a_1, \dots, a_n)$ интерполяционное приближение.

**) Элементы x_{ij} таблицы зависят не только от номера n , что подразумевается, но для простоты не отмечено в (4.1.5). С изменением n могут изменяться как сами элементы, так и длины строк. Некоторые строки могут оказаться пустыми. Таблицу можно было бы обозначить X_{k_0, k_1, \dots, k_m} , в тексте же принято более краткое обозначение X_n и число $n = k_0 + k_1 + \dots + k_m$, равное количеству исходных значений функции и производных, принято за номер таблицы.

$$X_n = \begin{pmatrix} x_{01} & x_{02} & \dots & x_{0k_0} \\ x_{11} & x_{12} & \dots & x_{1k_1} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mk_m} \end{pmatrix} \quad (4.1.5)$$

и семейством функций $\varphi_n(x, a_1, a_2, \dots, a_n)$.

Будем говорить, что задан интерполяционный процесс, если дана последовательность таблиц X_n для всяких $n=1, 2, \dots$ (или для избранной последовательности значений n) и последовательность соответствующих им семейств функций $\varphi_n(x, a_1, \dots, a_n)$, обладающих указанными выше свойствами.

За множество Φ принимается теоретико-множественная сумма семейств функций (4.1.3), что условно можно записать в виде

$$\Phi = \bigcup_n \varphi_n(x, a_1, \dots, a_n).$$

Правило выбора на шаге номера n из Φ интерполяционного приближения φ описано выше: из Φ выделяется семейство $\varphi_n(x, a_1, \dots, a_n)$, для него по таблице X_n составляется система (4.1.4), из которой находятся параметры a_1, \dots, a_n . Пусть это будут значения $a_1^{(n)}, \dots, a_n^{(n)}$. Для получения интерполяционного приближения осталось найденные значения подставить в уравнение семейства (4.1.3):

$$\varphi^{(n)} = \varphi_n(x, a_1^{(n)}, \dots, a_n^{(n)}).$$

Чтобы оценить близость интерполяционного приближения $\varphi^{(n)}$ к f , нужно, как обычно, ввести меру погрешности приближения. Если мы хотим выполнить интерполирование f в одной точке x , то естественной мерой погрешности является абсолютная величина разности:

$$|f(x) - \varphi^{(n)}(x)| = \rho(f, \varphi^{(n)}).$$

Более многообразно может быть определена мера погрешности в случае интерполяционного приближения на отрезке $\langle a, b \rangle$. Будем считать сейчас, что нами избрана какая-либо мера погрешности $\rho(f, \varphi^{(n)})$, и отметим, что в последующем для конечного отрезка $[a, b]$ мы будем пользоваться почти исключительно метрикой C , полагая

$$\rho(f, \varphi^{(n)}) = \max_x |f(x) - \varphi^{(n)}(x)|.$$

Примем следующее обычное определение.

Множество Φ называется всюду плотным в F , если для любой функции $f \in F$ и любого $\varepsilon > 0$ существует такое n и такие значения a_1, a_2, \dots, a_n , что будет

$$\rho(f, \varphi_n(x, a_1, a_2, \dots, a_n)) < \varepsilon. \quad (4.1.6)$$

Иными словами говоря, для каждой $f \in F$ существует последовательность функций φ_n из Φ , которая сходится к f в принятой метрике.

Если семейства φ_n выбраны при интерполировании так, что множество Φ этим свойством не обладает, то трудно рассчитывать на то, что интерполяционный процесс будет сходящимся, т. е. будет $\rho(f, \varphi^{(n)}) \rightarrow 0$ ($n \rightarrow \infty$).

В самом деле, в F будут существовать такие функции f , для каждой из которых нельзя будет построить последовательности функций $\varphi = \varphi_n(x, a_1, \dots, a_n)$, сходящейся к f в оценке ρ , как бы мы ни выбирали n и постоянные a_1, \dots, a_n . Тогда никакая последовательность интерполяционных приближений $\varphi^{(n)}$ не может сходиться к f .

Но если Φ обладает свойством плотности всюду в F , то этого еще недостаточно, чтобы имела место сходимост интерполяционной последовательности.

В самом деле, пусть $f \in F$. Тогда существует последовательность $\varphi_n(x, a_1, a_2, \dots, a_n)$, сходящаяся к f . Каждая из φ_n отвечает некоторым значениям a_1, a_2, \dots, a_n . Возьмем какую-нибудь из $\varphi_n = \varphi_n(x, a_1, \dots, a_n)$ и рассмотрим интерполяционное приближение $\varphi^{(n)} = \varphi^{(n)}(x, a_1^{(n)}, \dots, a_n^{(n)})$

того же номера n . При построении $\varphi^{(n)}$ параметры $a_k^{(n)}$ берутся не произвольно, а находятся из системы (4.1.4). Они зависят от узлов x_{ij} , а последние могут быть расположены так, что $a_k^{(n)}$ будут иметь значения, далекие от значений a_k , входящих в $\varphi_n(x, a_1, \dots, a_n)$. Поэтому интерполяционное приближение $\varphi^{(n)}$ может оказаться сильно отличающимся от $\varphi_n(x, a_1, \dots, a_n)$. Если взятая последовательность φ_n сходится к f , то последовательность соответствующих интерполяционных приближений не обязана обладать такой сходимостью.

Предшествующее изложение относилось к общему случаю, когда семейства функций (4.1.3) могли зависеть от параметров a_k любым образом. Картина построения приближения $\varphi^{(n)}$ значительно упрощается, если $\varphi_n(x, a_1, \dots, a_n)$ линейно зависит от a_k , так как тогда параметры находятся из линейной системы, и становится особенно простой и наглядной в следующем частном случае.

Пусть на $[a, b]$ дана последовательность линейно независимых m -кратно дифференцируемых функций $\omega_1(x), \omega_2(x), \dots, \omega_n(x), \dots$ *) За семейство (4.1.3) примем линейную комбинацию с произвольными коэффициентами первых n функций ω_k :

*) Функции ω_k бесконечной системы называются линейно независимыми, если любой конечный отрезок этой системы состоит из линейно независимых функций.

$$\varphi_n = a_1 \omega_1(x) + \dots + a_n \omega_n(x). \quad (4.1.7)$$

Система (4.1.4) для нахождения численных значений a_k будет состоять из n линейных уравнений

$$a_1 \omega_1^{(i)}(x_{ij}) + a_2 \omega_2^{(i)}(x_{ij}) + \dots + a_n \omega_n^{(i)}(x_{ij}) = f^{(i)}(x_{ij}) \quad (4.1.8)$$

и условием ее разрешимости при всяких $f^{(i)}(x_{ij})$ будет неравенство нулю определителя системы. Решение тогда будет единственным.

Множеством Φ здесь будет множество всевозможных конечных линейных комбинаций из ω_k вида (4.1.7).

Функции $\omega_k(x)$ должны быть выбраны так, чтобы линейные комбинации их были всюду плотны в F в принятой мере погрешности ρ .

Приведем примеры.

1. Рассмотрим систему степеней переменной x : $1, x, x^2, \dots, x^n, \dots$. Это — линейно независимые функции. Линейная комбинация первых n из них есть многочлен степени $n-1$:

$$a_0 x^{n-1} + a_1 x^{n-2} + \dots + a_{n-1} = P_{n-1}(x).$$

Интерполирование при помощи многочленов называется алгебраическим.

Множество Φ есть множество всех многочленов с действительными коэффициентами.

В математическом анализе доказывается теорема.*)

Если отрезок $[a, b]$ конечный и замкнутый и функция f непрерывна там вместе с производными первых m порядков, то для всякого $\varepsilon > 0$ существует многочлен некоторой степени n $P_n(x)$, для которого при всяких $x \in [a, b]$ выполняются неравенства

$$|f^{(i)}(x) - P_n^{(i)}(x)| < \varepsilon \quad (i=0, 1, \dots, m).$$

Эта теорема позволяет надеяться на то, что алгебраическое интерполирование может дать, по крайней мере в некоторых случаях, хорошее средство для вычисления не только значений самой функции, но и производных от нее всех порядков, когда они существуют и непрерывны.

2. В связи с интерполированием периодических функций, период которых мы считаем приведенным к 2π , рассмотрим систему тригонометрических функций $1, \cos x, \sin x, \cos 2x, \sin 2x, \dots$. Они линейно независимы. Линейная комбинация первых $2n+1$ их есть тригонометрический многочлен степени n :

*) Она является простым следствием известной теоремы о том, что на конечном отрезке $[a, b]$ можно равномерно и сколь угодно точно приблизиться при помощи многочлена ко всякой непрерывной на $[a, b]$ функции.

$$\varphi_{2n+1}(x) = a_0 + \sum_{k=1}^n (a_k \cos kx + b_k \sin kx) = p_n(x).$$

Интерполирование при помощи $p_n(x)$ называется тригонометрическим.

В этом случае Φ есть множество всех тригонометрических многочленов. Для них верна теорема о приближении при помощи $p_n(x)$ непрерывно дифференцируемых периодических функций, аналогичная той, которая приведена в примере 1 для алгебраических многочленов. Она позволяет надеяться, что тригонометрическое интерполирование может оказаться полезным для вычисления значений периодических функций и их производных.

Выше было описано множество Φ , из которого выбираются интерполяционные приближения $\varphi^{(n)}$, и указаны требования, обычно предъявляемые к Φ . Указаны также условия (4.1.4), при помощи которых находятся $\varphi^{(n)}$. Этим заканчивается лишь первая часть подготовительной работы для вычислений и исследований. После этого $\varphi^{(n)}(x)$ подвергают преобразованиям, объяснить которые в общем виде затруднительно и сделать это легче на частных типах интерполирования. Более подробно с такими преобразованиями можно ознакомиться на примере алгебраического интерполирования в следующих параграфах этой главы. Сейчас же мы вынуждены ограничиться несколькими, быть может недостаточно поучительными, общими фразами, объясняющими лишь цель таких преобразований.

В основе всего лежит тот факт, что не существует аналитического выражения для $\varphi^{(n)}(x)$, которое было бы удобным для всех случаев. Так как целей, для достижения которых употребляются приближения $\varphi^{(n)}(x)$, существует много, то это вынуждает строить многочисленные представления $\varphi^{(n)}(x)$, приспособленные к разным задачам: придавать $\varphi^{(n)}(x)$ форму, удобную для вычислений на быстродействующих машинах или вручную, удобную для оценки погрешности в наиболее важных классах функции, удобную для исследования сходимости интерполяционных процессов в этих классах и т. д. Все эти вопросы имеют, несомненно, технический характер, но они важны в вычислениях и исследованиях.

4.1.2. Остаток интерполирования. Сходимость интерполяционного процесса

Разница между функцией f и ее интерполяционным приближением $\varphi^{(n)}$ называется *остатком* или *погрешностью* интерполирования

$$f(x) - \varphi^{(n)}(x) = R_n(x).$$

При наших предположениях о f и Φ и ввиду условий (4.1.2) R есть m -кратно дифференцируемая функция на $\langle a, b \rangle$, удовлетворяющая условиям

$$R^{(i)}(x_{ij}) = 0 \quad (j=1, 2, \dots, k_i, i=0, 1, \dots, m).$$

Заметим прежде всего, что остаток является величиной весьма сложной природы, зависящей от многих факторов: от свойств функции f , от выбора семейства $\varphi_n(x, a_1, a_2, \dots, a_n)$, при помощи которого выполняется интерполирование, от таблицы X_n узлов x_{ij} , в частности от расположения x_{ij} на отрезке $\langle a, b \rangle$, что, как мы узнаем ниже, сильно влияет на значение и на поведение остатка как функции от x , и, наконец, от положения точки интерполирования x .

С проблемой изучения остатка интерполирования $R_n(x)$ тесно связана другая проблема, которая в значительной мере является ее следствием и в которой изучаются вопросы поведения остатка при изменении n , в частности вопросы сходимости интерполирования. Заключения о сходимости нередко вытекают из оценок остатка. Задача о сходимости является многосторонней, с различными выборами критерия сходимости. Чтобы показать, насколько разнообразно могут быть поставлены здесь исследования, мы приведем некоторые примеры.

Наиболее просто формулируется задача в том случае, когда мы стремимся найти только значения функции. Если x фиксировано и речь идет о нахождении значения f в этой точке, то $f(x)$, $\varphi^{(n)}(x)$, $R_n(x)$ будут численными величинами и мы должны определить, будет ли последовательность чисел $R_n(x) = f(x) - \varphi^{(n)}(x)$ сходиться к нулю. Первый вопрос, на который должна здесь ответить теория интерполирования, имеет следующее содержание: нужно определить, какими свойствами должны обладать функция f , последовательность семейств $\varphi_n(x, a_1, \dots, a_n)$ и таблиц X_n , чтобы было $\varphi^{(n)}(x) \rightarrow f(x)$.

Отметим попутно, что для вычислений доказательство сходимости равносильно доказательству принципиальной возможности найти $f(x)$ сколь угодно точно посредством избранного интерполяционного процесса при всяких n , больших некоторого числа.

Если в конкретной задаче или классе задач будет получен положительный или отрицательный ответ на вопрос о сходимости, то он, как правило, является только началом дальнейших исследований: нужно найти, как оценить скорость сходимости, каким будет асимптотическое представление остатка при больших n , какова точная оценка остатка и каким нужно взять n , чтобы найти $f(x)$ с заданной точностью, если сходимость медленная, то как ее можно ускорить, как можно целесообразно найти численное значение $f(x)$.

Если последовательность $\varphi^{(n)}(x)$ расходится, то в вопросе вычисления $f(x)$ еще не все потеряно, так как известно немало средств для нахождения у расходящейся последовательности «обобщенного предела», роль которого в нашей задаче играет $f(x)$.

Таков далеко не полный перечень вопросов, которые могут возникнуть после того, как решена основная задача о сходимости или расходимости $\varphi^{(n)}(x)$.

Все указанные выше вопросы являются общими и возникают при многих процессах приближений. В частности, они возникают и при других

задачах о сходимости приближений, о которых говорится ниже. Упомянули же о них мы в главе об интерполировании потому, что они особенно много и особенно успешно изучались в теории приближения и интерполирования функций.

Допустим теперь, что x может лежать в любом месте на $\langle a, b \rangle$ и речь идет о построении интерполяционного приближения к f на всем отрезке. Если такое приближение нам нужно для вычисления значений $f(x)$, то мы должны рассматривать поточечную сходимость $\varphi^{(n)}(x) \rightarrow f(x)$ на $\langle a, b \rangle$.

Если же мы ставим вопрос о равномерной сходимости, которая для вычислений имеет большее значение, то за меру погрешности может быть принята величина $\sup_x |R_n(x)| = \rho(f, \varphi^{(n)})$.

В некоторых приложениях интерес представляет не равномерная и даже не поточечная сходимость, а сходимость в смысле стремления к нулю среднего квадратичного отклонения $\varphi^{(n)}$ от f , что равносильно стремлению к нулю следующего интеграла:

$$\int_a^b [f(x) - \varphi^{(n)}(x)]^2 dx = \int_a^b R_n^2(x) dx.$$

Эту величину и принимают за меру оценки погрешности в таких задачах. При стремлении ее к нулю поточечной сходимости $\varphi^{(n)}(x) \rightarrow f(x)$ всюду на $\langle a, b \rangle$ здесь может не быть.

В научных и технических задачах иногда возникает потребность по таблице значений $f(x_k)$ ($k=1, 2, \dots, n$) вычислить не только значения $f(x)$ в нетабличных точках $\langle a, b \rangle$, но и значения производной $f'(x)$ во всех точках отрезка $\langle a, b \rangle$. Для этого часто по $f(x_k)$ составляют интерполяционное приближение $\varphi^{(n)}(x)$, значения которого принимают за значения f , а производную $\varphi^{(n)'}(x)$ — за $f'(x)$. За меру погрешности вычисления обеих функций f и f' , если мы заинтересованы в равномерном к ним приближении, может быть принята величина

$$\rho(f, \varphi^{(n)}) = p_0 \sup_x |f(x) - \varphi^{(n)}(x)| + p_1 \sup_x |f'(x) - \varphi^{(n)'}(x)|$$

$$(p_0, p_1 \geq 0, p_0 + p_1 = 1),$$

где p_0, p_1 — весовые коэффициенты.

Можно было бы привести еще примеры задач, в которых понятию сходимости придается разное содержание.

Возвратимся к проблеме изучения остатка $R_n(x)$. Его оценка является одной из основных задач. Выше мы обращали внимание на то, что остаток зависит от большого числа факторов. Остановим свое внимание на зависимости его от функции f . При получении оценки используются

свойства f , такие, как непрерывность, дифференцируемость, аналитичность и др. Поэтому каждая оценка рассчитана на класс функций, обладающих использованными в оценке свойствами. Чем шире класс функций, тем меньше, вообще говоря, будет существовать свойств, общих для всех функций класса, и тем грубее будет оценка $R_n(x)$ в этом классе, так как оценка должна учитывать функции с «самыми плохими свойствами». Такие оценки полезны в исследованиях сходимости интерполяционных процессов, так как мы при этом заинтересованы в доказательствах сходимости для возможно широких классов функций. Но такого рода оценки не могут принести большую пользу, когда мы при помощи их попытаемся определить, какое значение n нужно взять, чтобы погрешность $f(x) - \varphi^{(n)}(x) = R_n(x)$ была по абсолютной величине меньше заданной границы. Мы будем весьма часто получать завышенное значение номера шага n и рискуем проделать много излишней вычислительной работы для получения нужного результата.

В этой последней задаче определения n , так же как и при определении скорости сходимости интерполяционных процессов, полезными являются оценки остатка в более узких классах функций, важных в прикладном или принципиальном отношении. Примерами таких классов могут служить функции, аналитические на $\langle a, b \rangle$, функции, m -кратно непрерывно дифференцируемые там, абсолютно непрерывные и др. Ниже для некоторых типов интерполирования будут даны точные оценки остатков $R_n(x)$ в отдельных классах функций. Средством для их получения будет служить представление остатка, характерное для рассматриваемого класса функций, т. е. такое представление $R_n(x)$, которое верно для всех функций взятого класса и только для них.

Коротко остановимся на зависимости остатка интерполирования R_n от выбора семейства $\varphi_n(x, a_1, \dots, a_n)$.

До последних трех десятилетий усилия были направлены почти исключительно на исследование остатка алгебраического и, в значительно меньшей степени, тригонометрического интерполирования. Лишь в тридцатых и сороковых годах текущего века были получены достаточно общие представления $R_n(x)$, позволяющие судить о том, как будет изменяться форма остатка при изменении семейства $\varphi_n(x, a_1, \dots, a_n)$ и какой аппарат следует избрать для получения точных оценок остатка.*) Эти результаты найдены для того случая, когда семейство φ_n зависит от a_1, \dots, a_n линейно:

$$\varphi_n(x, a_1, \dots, a_n) = a_1 \omega_1(x) + a_2 \omega_2(x) + \dots + a_n \omega_n(x).$$

Последний вопрос, на котором мы остановим внимание во вступительном параграфе, — это зависимость погрешности интерполирования от выбора

*) Е. И. Ремез. О некоторых классах линейных функционалов в пространствах C_p и об остаточных членах формул приближенного анализа. Тр. ин-та математики АН УССР, 1939, т. 3, с. 21—62 и 1940, т. 4, с. 47—62.

узлов x_{ij} . Две следующие задачи имеют здесь, по-видимому, наибольший интерес: построение представлений остатка, дающих возможность достаточно наглядно судить о зависимости его от узлов, и такой выбор узлов, для которого можно ожидать наименьшего значения $R_n(x)$. Последняя задача обычно решается в смысле, который мы поясним на примере, где взята простейшая мера погрешности.

Пусть рассматривается класс F функций f и для каждой из них взято интерполяционное приближение $\varphi^{(n)}$ на отрезке $\langle a, b \rangle$. величиной, характеризующей точность интерполирования каждой функции f , будет $\sup_x |R_n(x)|$. За величину же, по которой судят о точности интерполирования всего семейства, принимают обычно $R = \sup_f \sup_x |R_n(x)|$. Она зависит от узлов x_{ij} , и их выбирают так, чтобы величина R имела наименьшее значение. Такой выбор узлов, как показали исследования частных случаев, является часто весьма целесообразным.

§ 4.2. КОНЕЧНЫЕ РАЗНОСТИ И РАЗНОСТНЫЕ ОТНОШЕНИЯ

В этом параграфе мы ознакомимся с некоторыми понятиями и терминами теории конечных разностей и разностных отношений. Для нас они будут иметь вспомогательное значение, поэтому мы ограничимся небольшим числом лишь самых необходимых сведений и изложение сделаем весьма кратким.

4.2.1. Конечные разности

Они применяются в исследованиях и при вычислении функций, заданных на сетке равноотстоящих точек. Как будет видно из приводимого ниже определения, конечные разности в вычислительной математике имеют значение, аналогичное дифференциалам в анализе бесконечно малых, и играют сходную роль.

Пусть дана сетка равноотстоящих точек с шагом $h > 0$ для аргумента x :

$$x_0, x_1 = x_0 + h, \dots, x_k = x_0 + kh, \dots$$

и известны соответствующие им значения функции $y = f(x)$:

$$y_0 = f(x_0), y_1 = f(x_0 + h), \dots, y_k = f(x_0 + kh), \dots$$

Конечными разностями первого порядка от функции $y = f(x)$ называются следующие величины:

$$\Delta y_0 = y_1 - y_0, \Delta y_1 = y_2 - y_1, \dots, \Delta y_k = y_{k+1} - y_k, \dots$$

Конечные разности от разностей первого порядка называются конечными разностями второго порядка

$$\Delta^2 y_0 = \Delta y_1 - \Delta y_0, \Delta^2 y_1 = \Delta y_2 - \Delta y_1, \dots, \Delta^2 y_k = \Delta y_{k+1} - \Delta y_k, \dots$$

и, вообще говоря, конечные разности от разностей порядка n называются конечными разностями порядка $n+1$

$$\Delta^{n+1} y_0 = \Delta^n y_1 - \Delta^n y_0, \Delta^{n+1} y_1 = \Delta^n y_2 - \Delta^n y_1, \dots, \Delta^{n+1} y_k = \Delta^n y_{k+1} - \Delta^n y_k, \dots$$

Укажем на некоторые легко проверяемые свойства конечных разностей.

1. Если $f(x) = u(x) + v(x)$, то для конечной разности $\Delta f(x) = f(x+h) - f(x)$ верно равенство

$$\Delta f(x) = \Delta u(x) + \Delta v(x).$$

2. Если $f(x) = Cu(x)$, где C — величина постоянная, то $\Delta f(x) = Cu(x)$.

Свойства 1 и 2 мы сформулировали для разностей первого порядка, но они, очевидно, верны для разностей любого порядка.

3. Если $y = P(x)$ есть многочлен степени n и x — величина произвольная, то конечная разность $\Delta P(x) = P(x+h) - P(x)$ есть многочлен от x степени $n-1$.

В силу свойств 1 и 2 утверждение достаточно проверить для степеней x , что делается весьма просто, так как

$$\Delta x^k = (x+h)^k - x^k = \frac{k}{1!} h x^{k-1} + \frac{k(k-1)}{2!} h^2 x^{k-2} + \dots$$

Отметим некоторые следствия, вытекающие отсюда. Если утверждение 3 применить дважды, то можно сказать, что разность второго порядка $\Delta^2 P(x)$ от многочлена степени n есть многочлен от x степени $n-2$ и т. д. Разность порядка n от многочлена степени n есть величина постоянная, и все разности, порядок которых больше n , будут равны нулю.

Выше было дано рекурсионное определение разностей всех порядков, но без труда может быть найдено их выражение непосредственно через значения функции. В самом деле,

$$\Delta y_0 = y_1 - y_0, \Delta^2 y_0 = \Delta y_1 - \Delta y_0 = (y_2 - y_1) - (y_1 - y_0) = y_2 - 2y_1 + y_0.$$

Выполняя несложную индукцию, убедимся в том, что для разности любого порядка верно следующее ее представление через значения функции:

$$\Delta^k y_0 = y_k - \frac{k}{1!} y_{k-1} + \frac{k(k-1)}{2!} y_{k-2} - \dots + (-1)^k y_0. \quad (4.2.1)$$

Его можно записать более компактно, при помощи «операции увеличения аргумента», а именно, введя операцию E , определенную равенством $Ef(x) = f(x+h)$. Произвольная действительная степень E^α операции

может быть определена равенством $E^\alpha f(x) = f(x + \alpha h)$. Применительно к значениям y_k функции эта операция дает $E y_k = y_{k+1}$ и $E^m y_k = y_{k+m}$.

Равенство (4.2.1) при помощи оператора E может быть коротко записано в следующей условной и простой форме:

$$\Delta^k y_0 = (E - 1)^k y_0. \quad (4.2.2)$$

Столь же просто может быть найдено выражение любого значения y_k функции через начальное ее значение y_0 и начальные значения конечных разностей $\Delta y_0, \Delta^2 y_0, \Delta^3 y_0, \dots$. В самом деле, по определению разности первого порядка $\Delta y_0 = y_1 - y_0$ имеем $y_1 = y_0 + \Delta y_0$. Далее, аналогично

$$\begin{aligned} y_2 &= y_1 + \Delta y_1 = (y_0 + \Delta y_0) + (\Delta y_0 + \Delta^2 y_0) = y_0 + 2\Delta y_0 + \Delta^2 y_0, \\ y_3 &= y_2 + \Delta y_2 = (y_0 + 2\Delta y_0 + \Delta^2 y_0) + (\Delta y_0 + 2\Delta^2 y_0 + \Delta^3 y_0) = \\ &= y_0 + 3\Delta y_0 + 3\Delta^2 y_0 + \Delta^3 y_0. \end{aligned}$$

Продолжив эти вычисления, по индукции найдем

$$y_k = y_0 + \frac{k}{1!} \Delta y_0 + \frac{k(k-1)}{2!} \Delta^2 y_0 + \dots + \Delta^k y_0 \quad (4.2.3)$$

или в условной форме

$$y_k = (1 + \Delta)^k y_0. \quad (4.2.4)$$

Некоторые сведения, касающиеся порядков малости конечных разностей и связи их с производными, будут приведены в конце параграфа, после выяснения аналогичных вопросов для разностных отношений.

4.2.2. Разностные отношения, их свойства и связь с конечными разностями

В том случае, когда значения аргумента являются не равноотстоящими, а произвольными, для исследования функций и вычислений вместо конечных разностей используют разностные отношения. Их часто называют также «разделенными разностями» и «подъемами» функций.

Пусть в произвольных попарно различных точках x_0, x_1, \dots известны значения функции f : $f(x_0), f(x_1), \dots$. Разностными отношениями первого порядка называются величины, имеющие смысл средних скоростей роста функции на соответствующих отрезках:

$$\begin{aligned} f(x_0, x_1) &= \frac{f(x_1) - f(x_0)}{x_1 - x_0}, \quad f(x_1, x_2) = \frac{f(x_2) - f(x_1)}{x_2 - x_1}, \\ f(x_2, x_3) &= \frac{f(x_3) - f(x_2)}{x_3 - x_2}, \dots \end{aligned}$$

По ним составляются разностные отношения второго порядка

$$f(x_0, x_1, x_2) = \frac{f(x_1, x_2) - f(x_0, x_1)}{x_2 - x_0},$$

$$f(x_1, x_2, x_3) = \frac{f(x_2, x_3) - f(x_1, x_2)}{x_3 - x_1}, \dots$$

Они, очевидно, связаны с изменением средней скорости роста функции при переходе от предыдущего отрезка (x_{i-1}, x_i) к следующему (x_i, x_{i+1}) . Более точно такую связь можно проследить по явному выражению разностных отношений через производные, которые будут даны ниже.

Разностные отношения третьего порядка определяются равенствами

$$f(x_0, x_1, x_2, x_3) = \frac{f(x_1, x_2, x_3) - f(x_0, x_1, x_2)}{x_3 - x_0},$$

$$f(x_1, x_2, x_3, x_4) = \frac{f(x_2, x_3, x_4) - f(x_1, x_2, x_3)}{x_4 - x_1}, \dots$$

Разностное отношение следующего порядка $k+1$ определяется через разностные отношения предыдущего порядка:

$$f(x_0, x_1, x_2, \dots, x_k, x_{k+1}) = \frac{f(x_1, x_2, \dots, x_{k+1}) - f(x_0, x_1, \dots, x_k)}{x_{k+1} - x_0}.$$

Укажем на некоторые свойства разностных отношений. Первые два из них являются очевидными, и мы сформулируем их без доказательства.

1. Свойство аддитивности. Если $f(x) = u(x) + v(x)$, то

$$f(x_0, x_1) = u(x_0, x_1) + v(x_0, x_1).$$

2. Свойство подобия. Если $f(x) = Cu(x)$, где C есть постоянная величина, то

$$f(x_0, x_1) = Cu(x_0, x_1).$$

Свойства 1 и 2 сформулированы для разностных отношений первого порядка, но они верны для разностных отношений любых порядков.

Третье свойство мы получим как следствие из представления разностного отношения через значения функции. Для разностного отношения первого порядка, по определению,

$$f(x_0, x_1) = \frac{f(x_1) - f(x_0)}{x_1 - x_0}.$$

Для разностного отношения второго порядка

$$\begin{aligned} f(x_0, x_1, x_2) &= \frac{f(x_1, x_2) - f(x_0, x_1)}{x_2 - x_0} = \\ &= \frac{1}{x_2 - x_0} \left\{ \left[\frac{f(x_1)}{x_1 - x_2} + \frac{f(x_2)}{x_2 - x_1} \right] - \left[\frac{f(x_0)}{x_0 - x_1} + \frac{f(x_1)}{x_1 - x_0} \right] \right\} = \\ &= \frac{f(x_0)}{(x_0 - x_1)(x_0 - x_2)} + \frac{f(x_1)}{(x_1 - x_0)(x_1 - x_2)} + \frac{f(x_2)}{(x_2 - x_0)(x_2 - x_1)}. \end{aligned}$$

При помощи индукции можно показать, что для любого k верно следующее представление:

$$\begin{aligned} f(x_0, x_1, \dots, x_k) &= \sum_{i=0}^k \frac{f(x_i)}{(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_k)} = \\ &= \sum_{i=0}^k \frac{f(x_i)}{\omega'(x_i)}, \quad \omega(x) = (x - x_0)(x - x_1) \dots (x - x_k). \end{aligned} \quad (4.2.5)$$

Если мы выполним любую перестановку аргументов x_0, x_1, \dots, x_k , то в последней части (4.2.5) переменятся местами слагаемые, что не изменит сумму. Это дает возможность сформулировать следующее свойство.

3. С в о й с т в о с и м м е т р и и. Разностное отношение $f(x_0, x_1, \dots, x_k)$ есть симметрическая функция аргументов x_0, x_1, \dots, x_k .

4. Если $f(x)$ есть многочлен степени n , то разностное отношение n -го порядка $f(x_0, x_1, \dots, x_n)$ не зависит от x_0, x_1, \dots, x_n и равняется коэффициенту при старшей степени x в многочлене f . Все разностные отношения порядка, большего n , равны нулю.

Это свойство может быть без труда доказано вычислениями, но мы его получим как простое следствие доказываемых ниже теорем (см. (4.2.6)).

Теорема 1. Если узлы x_0, x_1, \dots, x_k лежат на отрезке $[a, b]$ и $f(x)$ имеет непрерывную производную порядка k на $[a, b]$, то верно следующее представление разностного отношения порядка k через производную порядка k от f :

$$f(x_0, x_1, \dots, x_k) = \int_0^1 dt_1 \int_0^{t_1} dt_2 \dots \int_0^{t_{k-1}} dt_k f^{(k)} \left[x_0 + \sum_{i=1}^k t_i (x_i - x_{i-1}) \right]. \quad (4.2.6)$$

Заметим, что интеграл, стоящий справа, имеет смысл, так как область интегрирования есть k -мерная пирамида, определяемая неравенствами $0 \leq t_k \leq t_{k-1} \leq \dots \leq t_1 \leq 1$, и аргумент

$$\begin{aligned} x &= x_0 + \sum_{i=1}^k t_i (x_i - x_{i-1}) = (1-t_1)x_0 + (t_1-t_2)x_1 + \dots + (t_{k-1}-t_k)x_{k-1} + t_k x_k = \\ &= \frac{(1-t_1)x_0 + (t_1-t_2)x_1 + \dots + (t_{k-1}-t_k)x_{k-1} + t_k x_k}{(1-t_1) + (t_1-t_2) + \dots + (t_{k-1}-t_k) + t_k} \end{aligned}$$

равен среднему взвешенному значению, составленному из x_0, x_1, \dots, x_k с неотрицательными коэффициентами, и лежит на отрезке $[a, b]$, ввиду того что все x_i ($i=0, 1, \dots, k$) лежат на $[a, b]$.

Индуктивное доказательство (4.2.6) не имеет принципиальных трудностей и сложно лишь по записи. Мы ограничимся тем, что проверим равенство (4.2.6) для $k=1, 2$.

При $k=1$ правая часть (4.2.6) будет

$$\begin{aligned} \int_0^1 dt f'[x_0 + t(x_1 - x_0)] &= \left| \frac{1}{x_1 - x_0} f[x_0 + t(x_1 - x_0)] \right|_0^1 = \\ &= \frac{f(x_1) - f(x_0)}{x_1 - x_0} = f(x_0, x_1) \end{aligned}$$

и равенство верно.

При $k=2$, если интегрировать сначала по t_2 , будет

$$\begin{aligned} \int_0^1 dt_1 \int_0^{t_1} dt_2 f''[x_0 + t_1(x_1 - x_0) + t_2(x_2 - x_1)] &= \int_0^1 dt_1 \left| \frac{1}{x_2 - x_1} f'[x_0 + \right. \\ &+ t_1(x_1 - x_0) + t_2(x_2 - x_1)] \Big|_0^{t_1} = \frac{1}{x_2 - x_1} \left\{ \int_0^1 dt_1 f'[x_0 + t_1(x_2 - x_0)] - \right. \\ &- \left. \int_0^1 dt_1 f'[x_0 + t_1(x_1 - x_0)] \right\} = \frac{f(x_0, x_2) - f(x_1, x_0)}{x_2 - x_1} = \\ &= f(x_1, x_0, x_2) = f(x_0, x_1, x_2). \end{aligned}$$

Для $k=2$ равенство также верно.

Теорема 2. Если выполнены условия теоремы 1, то на $[a, b]$ существует такая точка ξ , что для разностного отношения порядка k верно равенство

$$f(x_0, x_1, \dots, x_k) = \frac{1}{k!} f^{(k)}(\xi). \quad (4.2.7)$$

Доказательство. Применим к интегралу (4.2.6) теорему о среднем значении. Он будет равен значению $f^{(k)}$ в некоторой средней точке области, умноженному на интеграл от единицы. Но для любой точки (t_1, t_2, \dots, t_k) области аргумент $x_0 + \sum_{i=1}^k t_i(x_i - x_{i-1})$ принадлежит отрезку $[a, b]$. Поэтому на $[a, b]$ существует такая точка ξ , что

$$f(x_0, x_1, \dots, x_k) = f^{(k)}(\xi) \int_0^1 dt_1 \int_0^{t_1} dt_2 \dots \int_0^{t_{k-1}} dt_k = \frac{1}{k!} f^{(k)}(\xi).$$

Отсюда сразу же следует свойство 4: если $f(x) = a_0 x^n + a_1 x^{n-1} + \dots$, то $f^{(n)}(x) = n! a_0$ и $f(x_0, x_1, \dots, x_n) = \frac{1}{n!} n! a_0 = a_0$.

Приведем теперь выражение произвольного значения функции $f(x_k)$ через начальное ее значение $f(x_0)$ и начальные значения разностных отношений $f(x_0, x_1)$, $f(x_0, x_1, x_2)$, $f(x_0, x_1, x_2, x_3)$, \dots . По определению, $f(x_0, x_1) = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$ и, следовательно, $f(x_1) = f(x_0) + (x_1 - x_0)f(x_0, x_1)$.

Ввиду полученного результата и определения $f(x_0, x_1, x_2)$ можно написать

$$\begin{aligned} f(x_2) &= f(x_1) + (x_2 - x_1)f(x_1, x_2) = [f(x_0) + (x_1 - x_0)f(x_0, x_1)] + \\ &+ (x_2 - x_1)[f(x_0, x_1) + (x_2 - x_0)f(x_0, x_1, x_2)] = f(x_0) + (x_2 - x_0)f(x_0, x_1) + \\ &+ (x_2 - x_0)(x_2 - x_1)f(x_0, x_1, x_2) \text{ и т. д.} \end{aligned}$$

При помощи несложных индуктивных рассуждений можно показать, что при всяком k будет

$$\begin{aligned} f(x_k) &= f(x_0) + (x_k - x_0)f(x_0, x_1) + (x_k - x_0)(x_k - x_1)f(x_0, x_1, x_2) + \\ &+ \dots + (x_k - x_0)(x_k - x_1) \dots (x_k - x_{k-1})f(x_0, x_1, \dots, x_k). \end{aligned} \quad (4.2.8)$$

Установим, наконец, связь между разностными отношениями и конечными разностями. Предположим, что значения аргумента $x_0, x_1, \dots, x_k, \dots$ являются равноотстоящими: $x_0, x_1 = x_0 + h, \dots, x_k = x_0 + kh, \dots$. Тогда

$$\begin{aligned} f(x_0, x_0 + h) &= \frac{f(x_0 + h) - f(x_0)}{(x_0 + h) - x_0} = \frac{y_1 - y_0}{h} = \frac{\Delta y_0}{1! h}, \\ f(x_0, x_0 + h, x_0 + 2h) &= \frac{f(x_0 + h, x_0 + 2h) - f(x_0, x_0 + h)}{(x_0 + 2h) - x_0} = \\ &= \frac{\Delta y_1 - \Delta y_0}{2h \cdot 1! h} = \frac{\Delta^2 y_0}{2! h^2} \dots \end{aligned}$$

и при любом k верно равенство

$$f(x_0, x_0 + h, \dots, x_0 + kh) = \frac{\Delta^k y_0}{k! h^k} \quad (y_i = f(x_0 + ih)). \quad (4.2.9)$$

Из теорем 1 и 2 для разностного отношения и (4.2.9) для конечной разности вытекает формулируемая ниже теорема о связи между $\Delta^k y_0$ и производной порядка k .

Теорема 3. Если f имеет на отрезке $[x_0, x_0 + kh]$ непрерывную производную порядка k , то для конечной разности порядка k верно следующее выражение ее через производную порядка k :

$$\Delta^k y_0 = k! h^k \int_0^1 dt_1 \int_0^{t_1} dt_2 \dots \int_0^{t_{k-1}} dt_k f^{(k)} \left[x_0 + h \sum_{i=1}^k t_i \right] \quad (4.2.10)$$

и на $[x_0, x_0 + kh]$ существует точка ξ , такая, что верно равенство

$$\Delta^k y_0 = h^k f^{(k)}(\xi). \quad (4.2.11)$$

По поводу полученного результата полезно сделать замечание. Если h есть малая величина, то конечные разности от функции f порядка k будут малыми величинами и представляет интерес выяснение закона их изменения при убывании h . Равенство (4.2.10) показывает, что если $f^{(k)}$ не обращается в нуль в точке x_0 , то $\Delta^k y_0$ будет малой величиной порядка k сравнительно с h , и если $f^{(k)}(x_0) = 0$, то $\Delta^k y_0$ будет иметь порядок малости, больший чем k .

§ 4.3. АЛГЕБРАИЧЕСКОЕ ИНТЕРПОЛИРОВАНИЕ ПО ЗНАЧЕНИЯМ ФУНКЦИИ. ПОГРЕШНОСТЬ ИНТЕРПОЛИРОВАНИЯ

4.3.1. Введение

В этом и нескольких следующих параграфах будут рассматриваться простейшие задачи интерполирования функций при помощи алгебраических многочленов. Начнем с проблемы интерполирования функции по нескольким ее значениям.

Допустим, что на конечном отрезке $[a, b]$ рассматривается функция $f(x)$, для которой либо строится приближение на всем отрезке, либо вычисляются с заданной точностью ее значения в нескольких точках отрезка.

Возьмем $n+1$ попарно разных узлов x_0, x_1, \dots, x_n на $[a, b]$ и будем считать известными соответствующие им значения функции $y_0=f(x_0)$, $y_1=f(x_1), \dots, y_n=f(x_n)$. Рассмотрим алгебраический многочлен степени n

$$P(x) = a_0x^n + a_1x^{n-1} + \dots + a_n, \quad (4.3.1)$$

содержащий $n+1$ неопределенных коэффициентов a_i . Выбор a_i подчиним требованиям совпадения значений P и f в узлах x_i :

$$P(x_k) = f(x_k) \quad (k=0, 1, \dots, n). \quad (4.3.2)$$

Они дадут для a_i систему $n+1$ линейных уравнений

$$\left. \begin{aligned} a_0x_0^n + a_1x_0^{n-1} + \dots + a_n &= f(x_0), \\ a_0x_1^n + a_1x_1^{n-1} + \dots + a_n &= f(x_1), \\ &\vdots \\ a_0x_n^n + a_1x_n^{n-1} + \dots + a_n &= f(x_n). \end{aligned} \right\} \quad (4.3.3)$$

Определитель системы

$$D = \begin{vmatrix} x_0^n & \dots & 1 \\ \vdots & \ddots & \vdots \\ x_n^n & \dots & 1 \end{vmatrix} = W_{n+1}(x_0, x_1, \dots, x_n)$$

есть определитель Вандермонда, и так как x_k различны между собой, он отличен от нуля и система имеет единственное решение.

Отсюда следует, что интерполирующий многочлен (4.3.1) может быть построен при любых узлах x_k , лишь бы они были различны, для любой функции f с конечными значениями в точках x_k и будет единственным.

Легко получить явное выражение для $P(x)$ через $x, x_k, f(x_k)$ при помощи определителей. Для этого присоединим (4.3.1) к системе (4.3.3) и полученные равенства

$$\begin{array}{l} -P(x) + a_0x^n + a_1x^{n-1} + \dots + a_n = 0, \\ -f(x_0) + a_0x_0^n + a_1x_0^{n-1} + \dots + a_n = 0, \\ .\quad.\quad.\quad.\quad.\quad.\quad.\quad.\quad.\quad.\quad.\quad.\quad.,\\ -f(x_n) + a_0x_n^n + a_1x_n^{n-1} + \dots + a_n = 0 \end{array}$$

будем рассматривать как однородную систему с неизвестными $-1, a_0, a_1, \dots, a_n$. Система заведомо имеет ненулевое решение, ввиду того что одно из неизвестных есть -1 , и поэтому ее определитель должен быть равен нулю. Если выписать этот определитель, приравнять его нулю и из полученного равенства найти $P(x)$, получим

$$P(x) = -\frac{1}{D} \begin{vmatrix} 0 & x^n & x^{n-1} & \dots & 1 \\ f(x_0) & \hline f(x_1) & \hline \dots & \hline f(x_n) & \hline \end{vmatrix}. \quad (4.3.4)$$

Такое представление $P(x)$ является сложным и требует вычисления определителей. Его мало применяют как в вычислениях, так и в теоретических исследованиях. Можно построить много других более простых представлений $P(x)$, удобных в разных отношениях. С некоторыми из них мы ознакомимся в п. 4.3.2 и в § 4.4.

Поясним теперь причины, побудившие избрать алгебраическое интерполирование в качестве основного объекта изучения среди других методов. Главным мотивом к этому был, несомненно, длительный, более чем двухвековой успешный опыт его применения к вычислениям при помощи таблиц, при составлении таблиц, при физических и инженерных расчетах и многом другом.

В основе большинства первых применений алгебраического интерполирования лежал простой факт, который легко можно понять, если воспользоваться формулой Тейлора. Пусть функция f — аналитическая или обладающая достаточно высоким порядком дифференцируемости, и нужно вычислить ее значения вблизи какой-либо точки, например x_0 . Она представима там в форме суммы:

$$f(x) = f(x_0) + (x-x_0)f'(x_0) + \frac{1}{2}(x-x_0)^2f''(x_0) + \frac{1}{6}(x-x_0)^3f'''(x_0) + \dots$$

Если x достаточно близок x_0 и разность $x - x_0$ настолько мала, что можно при принятой точности пренебречь всеми членами правой части, начиная со второго порядка малости, то $f(x)$ с нужной точностью совпадет с линейной функцией $f(x_0) + (x - x_0)f'(x_0)$ и для вычисления $f(x)$ достаточно будет выполнить интерполирование первой степени по двум ее

значениям $f(x_0)$ и $f(x_1)$, когда узел x_1 достаточно близок к x_0 . Если мы не можем пренебречь членом $\frac{1}{2}(x-x_0)^2 f''(x_0)$, но сможем отбросить все члены, содержащие $x-x_0$ в третьей и более высоких степенях, так что с нужной для нас точностью можно положить

$$f(x) \approx f(x_0) + (x-x_0)f'(x_0) + \frac{1}{2}(x-x_0)^2 f''(x_0),$$

то для вычисления $f(x)$ можно выполнить интерполирование второй степени при помощи трех значений $f(x_0)$, $f(x_1)$, $f(x_2)$ при условии достаточной близости x_1 и x_2 к x_0 и т. д.

Из проделанных несложных рассуждений вытекает, что, для того чтобы иметь возможность вычислить аналитическую или достаточно гладкую функцию в любой точке отрезка $[a, b]$, можно взять там достаточно густую сетку узлов x_0, x_1, \dots, x_n и пользоваться интерполированием разumno избранной степени. Чем более высокую степень интерполирования мы будем допускать, тем более редкую сетку точек можно будет взять.

Уже один простой факт возможности представления аналитических и гладких функций на всем отрезке таблицей их значений, часто небольшой по объему, должен был вызвать и вызвал большой интерес к алгебраическому интерполированию и побудил заняться разработкой его теории.

Как показало дальнейшее развитие математики, алгебраическое интерполирование имеет много большее значение, так как многочлены способны с любой заданной точностью представлять функции не только в «малых областях», а в сколь угодно больших, но конечных областях, при этом не только функции высокой гладкости, но и любые непрерывные, даже, может быть, не дифференцируемые. Соответствующая теорема указывалась нами в первом параграфе главы. В простой, но удобной для нас форме она может быть высказана в следующем виде.

Если f непрерывна на конечном замкнутом отрезке $[a, b]$ и ε есть любое положительное число, то существует такой многочлен $\Pi_m(x) = c_0 x^m + c_1 x^{m-1} + \dots + c_m$, для которого при всяких $x \in [a, b]$ выполняется неравенство $|f(x) - \Pi_m(x)| < \varepsilon$.

Но эта теорема оставляет открытым вопрос о том, можно ли достигнуть такого приближения при помощи интерполирования. Дело в том, что коэффициенты a_k интерполяционного многочлена нельзя задавать произвольно, так как они определяются узлами x_k ($k=0, 1, \dots, n$) и соответствующими им значениями функции $f(x_k)$ и должны быть найдены или из системы (4.3.3), или при помощи явного выражения (4.3.4). В интерполировании задача о равномерном приближении функции должна ставиться иначе, чем указано выше.

Пусть на замкнутом конечном отрезке $[a, b]$ задана непрерывная функция f и указано число $\varepsilon > 0$. Существует ли такое m и такие узлы x_0, x_1, \dots, x_m , что интерполяционный многочлен $P_m(x)$, построенный по этим узлам для функции f , будет выполнять неравенство $|f(x) - P_m(x)| < \varepsilon$ для всех $x \in [a, b]$?

Положительный ответ на поставленный вопрос, несомненно, повысил бы значение интерполирования в деле приближения непрерывных функций.

Если пользоваться только системой (4.3.3) или явным выражением (4.3.4) для $P(x)$, то найти ответ затруднительно. Но он легко получается, если воспользоваться двумя известными в конструктивной теории функций результатами.*)

Как будет показано в п. 4.8.4, многочлен наилучшего приближения $Q_m(x)$ степени m будет одновременно интерполирующим многочленом при некоторой системе $m+1$ узлов: $Q_m(x) = P_m(x)$.

Последовательность же многочленов $Q_m(x)$ равномерно сходится при $m \rightarrow \infty$ к $f(x)$ на $[a, b]$, и, следовательно, при всяком $\varepsilon > 0$ для всех достаточно больших m будет выполняться неравенство

$$|f(x) - Q_m(x)| = |f(x) - P_m(x)| < \varepsilon.$$

Этим доказано, что поставленный выше вопрос об интерполяционном приближении непрерывных функций имеет положительный ответ.

Чтобы правильно оценить такой результат, необходимо отметить, что он имеет в настоящее время только теоретическое значение, так как в нем нет никакого указания на то, как для заданной функции f эффективно находить соответствующие ей узлы x_k ($k=0, 1, \dots, n$), упомянутые в сформулированной выше проблеме. Устанавливается только существование таких узлов.

Более того, эти рассуждения заставляют думать, что может не существовать единого интерполяционного процесса, который мог бы обеспечить равномерную сходимость $P_n(x) \rightarrow f(x)$ на $[a, b]$ для всякой непрерывной функции. Забегая немного вперед, скажем, что такой процесс действительно невозможен. Этот вопрос и некоторые связанные с ним задачи будут изучаться в параграфе, посвященном проблеме сходимости интерполирования.

4.3.2. Интерполяционные формулы Лагранжа и Ньютона

Укажем сейчас два наиболее употребительных представления интерполяционного многочлена. Начнем с формулы Лагранжа. Введем сначала многочлены влияния отдельных узлов. Их называют часто коэффициентами Лагранжа.

*) Необходимые сведения можно найти в добавлении III.

Рассмотрим узел x_k . Многочлен влияния этого узла $\omega_k(x)$ определяется условиями: 1) степень его равна n и 2) он обращается в единицу при $x=x_k$ и в нуль во всех прочих узлах x_i ($i \neq k$). Так как число узлов x_i ($i \neq k$) равно n , а это есть всевозможные корни многочлена $\omega_k(x)$, и если принять во внимание условие $\omega_k(x_k)=1$, то будет ясно, что

$$\omega_k(x) = \frac{(x-x_0) \dots (x-x_{k-1})(x-x_{k+1}) \dots (x-x_n)}{(x_k-x_0) \dots (x_k-x_{k-1})(x_k-x_{k+1}) \dots (x_k-x_n)}.$$

$\omega_k(x)$ можно записать в более простой форме, если ввести многочлен $\omega(x)$, для которого узлы x_0, \dots, x_n будут всевозможными простыми нулями: $\omega(x) = (x-x_0)(x-x_1) \dots (x-x_n)$. Очевидно,

$$\omega_k(x) = \frac{\omega(x)}{(x-x_k)\omega'(x_k)}.$$

Теперь легко проверить правильность следующего выражения интерполяционного многочлена (4.3.3) через $\omega_k(x)$ и значения функции $f(x_k)$ ($k=0, 1, \dots, n$):

$$\begin{aligned} P(x) &= \omega_0(x)f(x_0) + \omega_1(x)f(x_1) + \dots + \omega_n(x)f(x_n) = \\ &= \sum_{k=0}^n \frac{\omega(x)}{(x-x_k)\omega'(x_k)} f(x_k). \end{aligned} \quad (4.3.5)$$

Действительно, каждый многочлен $\omega_k(x)$ имеет степень n , поэтому правая часть равенства есть многочлен степени не выше n . При $x=x_0$ будет $\omega_0(x_0)=1$, $\omega_k(x_0)=0$ ($k>0$) и правая часть будет равна $1 \cdot f(x_0) + 0 \cdot f(x_1) + \dots + 0 \cdot f(x_n) = f(x_0)$. Но узел x_0 по своему значению ничем не отличается от прочих узлов и, следовательно, при $x=x_k$ правая часть равна $f(x_k)$ ($k=0, 1, \dots, n$). Таким образом, правая часть удовлетворяет условиям (4.3.2). Этим (4.3.5) доказано.

Укажем на некоторые особенности лагранжевой формулы (4.3.5). Свойства интерполяционного многочлена зависят, очевидно, от двух факторов: от выбора узлов x_k и от интерполируемой функции f . В формуле (4.3.5) оба фактора разделены, так как многочлены $\omega_k(x)$ зависят только от узлов, а свойства функции f учитываются множителями $f(x_k)$. Это обстоятельство оказывается полезным в некоторых вопросах теории сходимости интерполирования, и формулой (4.3.5) там широко пользуются.

В отношении вычислений формула Лагранжа удобна в задаче интерполирования многих функций в одной точке x , так как значения множителей $\omega_k(x)$ можно вычислить однажды для всех функций. Но вычислительное применение (4.3.5) имеет существенный недостаток, так как нужно заранее определять число $n+1$ узлов, необходимое для достижения

принятой точности. Желание избежать ненужной затраты труда побуждает вычислителя стараться обойтись наименьшим числом узлов и нередко оказывается, что заданное им число узлов является недостаточным или бывает необходимо проверить точность полученного результата. В обоих случаях к взятым узлам добавляют еще один или несколько узлов и выполняют вычисления заново. Тогда в формуле (4.3.5) не только добавятся новые члены, но потребуется перевычислить все ранее найденные члены суммы, так как в них появятся новые множители.

Другое представление $P(x)$, к которому мы перейдем, допускает последовательное уточнение результатов вычислений и часто не требует при применениях предварительного указания степени интерполирования. По строению оно аналогично формуле Тейлора и обращается в нее, если перейти к пределу, когда все узлы интерполирования x_0, x_1, \dots, x_n будут стремиться к какому-либо одному значению, например x_0 .

Покажем, что интерполирующий многочлен можно записать в приводимом ниже виде, который называется формулой Ньютона:

$$P(x) = f(x_0) + (x-x_0)f(x_0, x_1) + (x-x_0)(x-x_1)f(x_0, x_1, x_2) + \\ + \dots + (x-x_0)(x-x_1)\dots(x-x_{n-1})f(x_0, x_1, \dots, x_n). \quad (4.3.6)$$

В правильности равенства проще всего убедиться путем его проверки. Достаточно показать, что многочлен $P(x)$, определенный равенством (4.3.6) и имеющий, очевидно, степень не большую n , удовлетворяет условиям (4.3.2). При $x=x_0$ все члены правой части, начиная со второго, обращаются в нуль и остается $P(x_0)=f(x_0)$. При $x=x_1$ справа останется $f(x_0) + (x_1-x_0)f(x_0, x_1)$, что, на основании (4.2.8) при $k=1$, равно $f(x_1)$. Значит, $P(x_1)=f(x_1)$. Продолжая такие вычисления и пользуясь (4.2.8), убедимся в том, что для любого $k=0, 1, \dots, n$ будет $P(x_k)=f(x_k)$.

Формула Ньютона имеет строение более сложное, чем (4.3.5), и требует составления разностных отношений $f(x_0, x_1, \dots, x_k)$ ($k=1, 2, \dots, n$). При добавлении к x_0, x_1, \dots, x_n нового узла x_{n+1} все ранее найденные члены сохраняются и в формуле добавляется еще один член

$$(x-x_0)(x-x_1)\dots(x-x_n)f(x_0, x_1, \dots, x_n, x_{n+1}).$$

Это позволяет не задавать заранее число узлов и постепенно увеличивать точность результата, добавляя последовательно по одному новому узлу.

4.3.3. Остаток интерполирования и его представления для некоторых классов функций

Свойства остатка, или погрешности интерполирования, $R(x)=f(x)-P(x)$ зависят от свойств функции f и от выбора узлов x_k ($k=0, 1, \dots, n$) и было бы желательно найти такие представления $R(x)$, которые учитывали бы заранее некоторые наиболее распространенные структурные свойства функций и позволяли бы без больших затруднений судить

о влиянии на остаток расположения на оси узлов x_k и точки x . Начнем с представления $R(x)$ для общего случая, не налагающего, по сути дела, никаких ограничений на f и предполагающего лишь то, что в точках x_0, x_1, \dots, x_n, x она имеет конечные значения. Такое представление ввиду его большой общности редко применяется в исследованиях. Для нас оно будет полезным в двух отношениях: во-первых, как источник получения специализированных представлений для более узких классов функций и, во-вторых, при его помощи легко может быть получено приближенное выражение для остатка, в котором все величины, входящие в него, являются вычислимыми.

С достаточной простотой нужная нам формула может быть получена, если воспользоваться равенством (4.2.8), дающим значение f в любой точке x_k через начальное значение функции $f(x_0)$ и разностные отношения $f(x_0, x_1), f(x_0, x_1, x_2), \dots, f(x_0, x_1, \dots, x_k)$. Применив его к точкам $x_0, x_1, x_2, \dots, x_n, x$, найдем

$$\begin{aligned} f(x) = & f(x_0) + (x-x_0)f(x_0, x_1) + (x-x_0)(x-x_1)f(x_0, x_1, x_2) + \dots + \\ & + (x-x_0) \dots (x-x_{n-1})f(x_0, x_1, \dots, x_n) + \\ & + (x-x_0) \dots (x-x_{n-1})(x-x_n)f(x_0, x_1, \dots, x_n, x). \end{aligned}$$

Если сравнить это равенство с (4.3.6), то будет видно, что сумма всех членов правой части равенства, кроме последнего, есть не что иное, как интерполирующий многочлен $P(x)$ в форме Ньютона. Поэтому последний член справа есть остаток $R(x)$. Это дает возможность высказать теорему об остатке $R(x)$.

Теорема 1. Если f есть любая функция с конечными значениями в точках x_0, x_1, \dots, x_n, x , то остаток $R(x)$ ее алгебраического интерполирования по значениям в точках x_k при помощи многочлена степени n представим в виде

$$\begin{aligned} R(x) = & (x-x_0)(x-x_1) \dots (x-x_n)f(x_0, x_1, \dots, x_n, x) = \\ = & \omega(x)f(x_0, x_1, \dots, x_n, x). \end{aligned} \quad (4.3.7)$$

Множитель $\omega(x)$ зависит только от узлов x_k ($k=0, 1, \dots, n$), со свойствами же функции f связано разностное отношение $f(x_0, x_1, \dots, x_n, x)$. Оно не может быть вычислено, так как зависит от $f(x)$, но для него, по крайней мере в некоторых случаях, может быть получено вычисляемое приближенное выражение, полезное как ориентировочное. Чтобы сделать более понятной наглядную сторону вопроса, мы проведем рассуждения в более ограничительных предположениях, чем выше, именно будем считать, что f имеет непрерывную производную порядка $n+1$ на отрезке, содержащем все узлы x_k и точку x , и применим к $f(x_0, x_1, \dots, x_n, x)$

соотношение (4.2.7), устанавливающее связь между разностным отношением и производной. В нашем случае оно дает

$$f(x_0, x_1, \dots, x_n, x) = \frac{1}{(n+1)!} f^{(n+1)}(\xi),$$

где ξ есть точка указанного выше отрезка. Ее положение зависит от x . Если на этом отрезке произвольная $f^{(n+1)}(x)$ мало изменяется, что наверное будет, когда x_k и x принадлежат малому отрезку, то $f^{(n+1)}(\xi)$ будет мало изменяться с изменением x и нужное нам разностное отношение может быть заменено другим, в котором вместо x может быть взято любое табличное значение x , например значение x_{n+1} : $f(x_0, x_1, \dots, x_n, x) \approx \approx f(x_0, x_1, \dots, x_n, x_{n+1})$. Это дает возможность для $R(x)$ указать приближенное выражение

$$R(x) \approx \omega(x) f(x_0, x_1, \dots, x_n, x_{n+1}), \quad (4.3.8)$$

значение которого может быть вычислено.

Получим теперь из (4.3.7) более специализированное представление $R(x)$, рассчитанное на функции высокого порядка дифференцируемости. Если говорить более точно, мы будем считать, что все x_k и x принадлежат отрезку $[c, d]$ и f имеет на $[c, d]$ непрерывную производную порядка $n+1$. К разностному отношению $f(x_0, x_1, \dots, x_n, x)$ может быть применена теорема 1 и формула (4.2.6), если в ней считать $k=n+1$ и $x_{n+1}=x$. Они дадут интегральное выражение для $f(x_0, x_1, \dots, x_n, x)$ через $f^{(n+1)}$ и позволят высказать теорему об остатке.

Теорема 2. Если точки x_0, x_1, \dots, x_n, x принадлежат отрезку $[c, d]$ и f имеет на $[c, d]$ непрерывную производную порядка $n+1$, остаток ее интерполирования при помощи многочлена степени n по узлам x_0, x_1, \dots, x_n представим в виде

$$R(x) = \omega(x) \int_0^1 dt_1 \int_0^{t_1} dt_2 \dots \int_0^{t_n} dt_{n+1} f^{(n+1)} \left[x_0 + \sum_{v=1}^{n+1} t_v (x_v - x_{v-1}) \right],$$

$$(x_{n+1} = x). \quad (4.3.9)$$

Полученный результат интересен в том отношении, что он не содержит никаких неизвестных величин и позволяет составить достаточно ясное представление о зависимости $R(x)$ от свойств производной $f^{(n+1)}(x)$.

Для $R(x)$ можно получить более простое выражение, если взять не интегральное представление $f(x_0, x_1, \dots, x_n, x)$, а воспользоваться теоремой 2 и более простой формулой (4.2.7), полагая в ней $k=n+1$ и $x_{n+1}=x$. Тогда получится

Теорема 3. Если выполняются условия теоремы 2, то на $[c, d]$ существует точка ξ , такая, что для остатка $R(x)$ интерполирования f при помощи многочлена степени n по узлам x_0, x_1, \dots, x_n верно равенство

$$R(x) = \frac{\omega(x)}{(n+1)!} f^{(n+1)}(\xi). \quad (4.3.10)$$

Равенство (4.3.10) называется формулой Лагранжа для остатка интерполирования. По сравнению с (4.3.9) она выигрывает в простоте, но проигрывает в точности информации, так как содержит величину ξ , о которой мы знаем лишь то, что она лежит на отрезке $[c, d]$.

Применим сейчас (4.3.10) к решению двух простых задач о выборе узлов интерполирования. Рассмотрим всевозможные функции, $n+1$ раз непрерывно дифференцируемые на $[a, b]$, с производной порядка $n+1$, ограниченной по модулю числом M : $|f^{(n+1)}(x)| \leq M$ ($x \in [a, b]$). В таком классе функций остаток интерполирования имеет оценку

$$|R(x)| \leq \frac{M}{(n+1)!} |x-x_0| \cdot |x-x_1| \dots |x-x_n|. \quad (4.3.11)$$

Она является точной и достигается в том случае, когда f есть многочлен степени $n+1$ вида

$$f(x) = \frac{M}{(n+1)!} x^{n+1} + c_1 x^n + c_2 x^{n-1} + \dots$$

Допустим теперь, что функция f дана нам таблицей значений, x есть не табличное значение аргумента и мы должны интерполировать f в точке x при помощи многочлена степени n , взяв за узлы x_0, x_1, \dots, x_n любые $n+1$ табличных значений аргумента. Как следует выбрать эти значения, чтобы погрешность интерполирования была наименьшей? Задача весьма простая, имеет очень наглядное значение, и ответ к ней может быть предсказан заранее: наименьшее значение $|R(x)|$ получится, вообще говоря, в случае, когда в качестве узлов x_0, x_1, \dots, x_n будут взяты $n+1$ табличных значений аргумента, ближайших к точке x .

Точные выражения остатка (4.3.9) и (4.3.10) сложно зависят от узлов x_k и малопригодны для получения из них правила выбора. Сосредоточим свое внимание на оценке (4.3.11). Так как оценка точная, то можно ожидать, что значения узлов, для которых она имеет наименьшую величину,

будут наилучшими при интерполировании. Множитель $\frac{M}{n+1!}$ не зависит

от выбора узлов, и мы должны x_k при фиксированном значении x выбрать так, чтобы произведение $|x-x_0| \cdot |x-x_1| \dots |x-x_n|$ имело наименьшее значение, когда в качестве x_k берутся табличные значения аргумента.

Это и подтверждает высказанное правило: за x_k следует взять табличные значения аргумента, ближайшие к x .

Рассмотренная простая задача представляет интерес не своим легко предсказываемым ответом, а скорее как испытание на правильность полученных представлений остатка и его оценки.

Вторая задача имеет не столь очевидный ответ и приводит к более интересным следствиям. Рассмотрим тот же класс функций, что и в первой задаче. Возьмем $n+1$ узлов x_0, x_1, \dots, x_n на $[a, b]$ и построим по ним интерполяционное приближение $P(x)$ для f на всем отрезке $[a, b]$. За меру точности приближения в точке x естественно принять $|R(x)| = |f(x) - P(x)|$, за меру же приближения на всем отрезке $[a, b]$ следует принять $\max_x |R(x)|$. Наконец, за величину, характеризующую погрешность приближения на $[a, b]$ всех функций семейства, должна быть взята величина

$$\sup_f \max_x |R(x)|. \quad (4.3.12)$$

Она зависит от выбора узлов x_0, x_1, \dots, x_n .

Поставим задачу: как следует выбрать узлы x_k ($k=0, 1, \dots, n$), чтобы их можно было признать наилучшими при построении интерполяционных приближений на $[a, b]$ всех функций f взятого класса. Такими узлами, очевидно, нужно признать те, при которых величина (4.3.12) имеет наименьшее значение.

Найдем (4.3.12). Из оценки (4.3.11), которую мы запишем коротко

$$|R(x)| \leq \frac{M}{(n+1)!} |\omega(x)|,$$

следует:

$$\max_x |R(x)| \leq \frac{M}{(n+1)!} \max_x |\omega(x)|.$$

Отметим, что неравенство переходит в равенство для случая, когда f есть указанный выше многочлен

$$f(x) = \frac{M}{(n+1)!} x^{n+1} + c_1 x^n + \dots$$

Так как правая часть неравенства не зависит от f , то

$$\sup_f \max_x |R(x)| \leq \frac{M}{(n+1)!} \max_x |\omega(x)|,$$

и так как оценка достижима, то следует взять точное равенство

$$\sup_f \max_x |R(x)| = \frac{M}{(n+1)!} \max_x |\omega(x)|.$$

От выбора узлов x_k ($k=0, 1, \dots, n$) в правой части зависит только множитель $\max_x |\omega(x)| = \max_x |(x-x_0)(x-x_1)\dots(x-x_n)|$ и для него ставится задача: среди всех многочленов $\omega(x) = x^{n+1} + b_1x^n + \dots$, корни которых принадлежат отрезку $[a, b]$, нужно найти тот, для которого $\max_{x \in [a, b]} |\omega(x)|$ будет наименьшим.

Заметим, что оговорка относительно расположения корней x_k многочлена $\omega(x)$ вызвана тем, что функции f мы считаем определенными на $[a, b]$ и интерполировать их можем только по узлам, лежащим на этом отрезке. Можно было бы изменить задачу об $\omega(x)$ и поставить ее так: среди всех многочленов $\omega(x) = x^{n+1} + b_1x^n + b_2x^{n-1} + \dots$ с произвольными коэффициентами b_1, b_2, \dots нужно найти тот многочлен, для которого $\max_{a \leq x \leq b} |\omega(x)|$ будет наименьшим.

Решением ее является многочлен, наименее уклоняющийся от нуля на $[a, b]$. Для отрезка $-1 \leq t \leq 1$ таким многочленом будет многочлен Чебышева первого рода $(n+1)$ -й степени

$$\bar{T}_{n+1}(t) = \frac{1}{2^n} \cos[(n+1) \arccos t] = t^{n+1} - \dots$$

Корнями его являются числа

$$t_k = \cos \frac{2k+1}{2(n+1)} \pi \quad (k=0, 1, \dots, n).$$

Они лежат внутри отрезка $[-1, 1]$. От отрезка $[-1, 1]$ линейным преобразованием $x = \frac{a+b}{2} + \frac{b-a}{2} t$ можно перейти к $[a, b]$ и получить

многочлен, наименее уклоняющийся от нуля на $[a, b]$. Корни такого многочлена будут лежать внутри отрезка $[a, b]$. Поэтому обе задачи о нахождении многочлена $\omega(x)$ — с оговоркой и без оговорки о положении корней на $[a, b]$ — равносильны и имеют одно и то же решение.

Мы пришли к заключению, что наилучшими узлами для построения интерполяционного приближения на $[a, b]$ функций f , удовлетворяющих условию $|f^{(n+1)}(x)| \leq M$, являются корни многочлена степени $n+1$, наименее уклоняющегося от нуля на отрезке $[a, b]$. Заключение не зависит от значения M и остается верным для всяких функций, $(n+1)$ -кратно непрерывно дифференцируемых на $[a, b]$.

В последующем мы увидим, что интерполяционные процессы, в которых выполняется интерполирование по узлам, являющимся корнями многочленов Чебышева первого рода возрастающих степеней, обладают целым рядом замечательных свойств. В частности, такие интерполяционные процессы будут равномерно сходиться для широких множеств функций, например, как будет показано в § 4.5, для всякой абсолютно непрерывной функции.*)

Остановимся еще на представлении остатка интерполирования, характерном для аналитических функций. Рассмотрим плоскость комплексной переменной $z = x + iy$ и допустим, что в ней указана конечная замкнутая область D , ограниченная спрямляемым контуром l и содержащая внутри себя отрезок $[a, b]$ действительной оси.

Предположим теперь, что в D определена однозначная аналитическая функция $f(z)$, регулярная всюду в D , включая и ее контур l . Нас будет интересовать задача интерполирования f только на отрезке $[a, b]$ действительной оси, мы не будем выходить с этого отрезка в комплексную плоскость и оставим в стороне задачу об интерполировании f вне отрезка $[a, b]$.

Роль, которую при этом будет играть область D , мы выясним ниже.

Как и выше, возьмем $n+1$ различных узлов x_0, x_1, \dots, x_n на $[a, b]$ и пусть z есть точка интерполирования, отличная от них. По ним составим интерполирующий многочлен

$$P(z) = \sum_{k=0}^n \frac{\omega(z)}{(z-x_k)\omega'(x_k)} f(x_k), \quad \omega(z) = (z-x_0)\dots(z-x_n). \quad (4.3.13)$$

Он определен на всей плоскости z и может быть принят за интерполяционное приближение f всюду в D . Несколько позже мы будем его рассматривать только на $[a, b]$, а сейчас z будем считать любой внутренней точкой области D , отличной от узлов x_k . Рассмотрим погрешность интерполирования $R(z) = f(z) - P(z)$. Это есть аналитическая функция z , регулярная всюду в D . Укажем для нее представление контурным интегралом.

Теорема 4. Если выполнены указанные выше условия для f , D , l и z , то для остатка $R(z)$ верно равенство

$$R(z) = \frac{\omega(z)}{2\pi i} \int_l \frac{f(t)}{\omega(t)(t-z)} dt. \quad (4.3.14)$$

Доказательство. Чтобы убедиться в правильности (4.3.14), достаточно вычислить правую часть и показать, что она равна $f(z)$ —

*) Функция называется абсолютно непрерывной, если она представима в форме неопределенного интеграла от суммируемой по Лебегу функции.

— $P(z)$. Отбросим на время множитель $\omega(z)$ и рассмотрим оставшийся после этого контурный интеграл. Он равен сумме вычетов функции $\frac{f(t)}{\omega(t)(t-z)}$ в особых точках, лежащих внутри D . $f(z)$ не имеет там особенностей, и особыми точками будут нули знаменателя $t=x_k$ ($k=0, 1, \dots, n$) и $t=z$. Все они простые, и вычеты в этих точках находятся по известным правилам, которые даются в теории функций комплексной переменной:

$$\begin{aligned}\operatorname{Res} \left[\frac{f(t)}{\omega(t)(t-z)} \right]_{t=z} &= \frac{f(z)}{\omega(z)}, \\ \operatorname{Res} \left[\frac{f(t)}{\omega(t)(t-z)} \right]_{t=x_k} &= \frac{f(x_k)}{\omega'(x_k)(x_k-z)}.\end{aligned}$$

Поэтому

$$\frac{\omega(z)}{2\pi i} \int_l \frac{f(t)}{\omega(t)(t-z)} dt = f(z) - \sum_{k=0}^n \frac{\omega(z) \cdot f(x_k)}{(z-x_k)\omega'(x_k)} = f(z) - P(z) = R(z).$$

Попутно отметим, что при помощи формулы Коши для функции f :

$$f(z) = \frac{1}{2\pi i} \int_l \frac{f(t)}{t-z} dt$$

и равенства (4.3.14) для остатка легко получается представление многочлена $P(z)$ контурным интегралом

$$P(z) = \frac{1}{2\pi i} \int_l \frac{\omega(t) - \omega(z)}{\omega(t)} \frac{f(t)}{t-z} dt. \quad (4.3.15)$$

Теперь вернемся к нашей основной задаче изучения остатка в точке x отрезка $[a, b]$ действительной оси

$$R(x) = \frac{\omega(x)}{2\pi i} \int_l \frac{f(t)}{\omega(t)(t-x)} dt. \quad (4.3.16)$$

Сначала — несколько чисто качественных замечаний. Вдали от особых точек аналитическая функция изменяется весьма плавно. Чем шире будет область D и чем дальше от $[a, b]$ будет удалена ее граница l , тем более плавным будет поведение f на $[a, b]$ и тем меньшей погрешности интерполирования можно тогда ожидать. Если же иметь в виду интер-

поляционный процесс при $n \rightarrow \infty$, то тем более быстро он должен сходиться.

Рассмотрим несколько более подробно картину поведения остатка. В его представлении (4.3.16) от выбора узлов x_k и их числа $n+1$ зависит величина

$$\frac{\omega(x)}{\omega(t)} = \prod_{k=0}^n \frac{x-x_k}{t-x_k}.$$

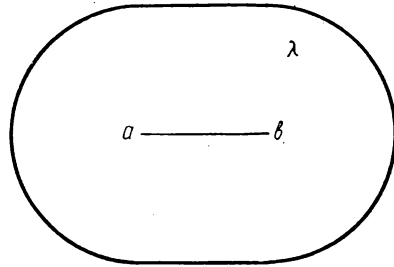


Рис. 4.3.1

Возьмем один сомножитель $\frac{x-x_k}{t-x_k}$. Здесь t — переменная точка контура l . Если контур l достаточно широкий, то при всяких k и всех положениях t на l отношение $\frac{x-x_k}{t-x_k}$ будет иметь малый модуль. Дробь $\frac{\omega(x)}{\omega(t)}$ будет весьма малой величиной, быстро стремящейся к нулю при росте n . Поэтому можно ожидать малой погрешности $R(x)$ интерполирования и быстрой сходимости интерполяционного процесса при $n \rightarrow \infty$.

Самая грубая оценка, которую здесь можно получить, — следующая. При всяком положении x на $[a, b]$ и любом k имеем $|x-x_k| \leq b-a$. Обозначим r расстояние от контура l до $[a, b]$. Всегда будет $|t-x_k| \geq r$ и

$$\left| \frac{x-x_k}{t-x_k} \right| \leq \frac{b-a}{r}.$$

Для $R(x)$ получится оценка

$$|R(x)| \leq \left(\frac{b-a}{r} \right)^{n+1} \frac{1}{2\pi} \int_l \frac{|f(t)|}{|t-x|} dt.$$

Когда $r > b-a$, то $R(x)$ будет при больших n малой величиной и будет стремиться к нулю не медленнее, чем убывает показательная функция, стоящая в неравенстве справа. Последнее означает, что если рассматривать интерполяционный процесс при $n \rightarrow \infty$ с любым выбором узлов x_k на $[a, b]$, то $P(x)$ будет равномерно на $[a, b]$ сходиться к $f(x)$.

Построим линию λ , все точки которой удалены от отрезка $[a, b]$ на расстояние $b-a$. Она состоит из двух полуокружностей радиусов $b-a$ с центрами в точках a и b и двух прямолинейных отрезков, параллельных $[a, b]$ и отстоящих от него на расстоянии $b-a$ (рис. 4.3.1).

Когда функция $f(z)$ является регулярной в замкнутой области, ограниченной линией λ , то для нее в качестве l может быть взята некоторая линия, охватывающая λ , и условие $r > b-a$ для l будет выполнено. Отсюда следует, что будет верно утверждение:

Если $f(z)$ регулярна в замкнутой области, ограниченной линией λ , то для нее интерполяционный процесс при $n \rightarrow \infty$ будет равномерно на $[a, b]$ сходиться к $f(x)$ при любом выборе узлов x_k ($k=0, 1, \dots, n$) на отрезке $[a, b]$.

Это утверждение мы привели здесь как самую простую иллюстрацию значения формулы (4.3.14) для остатка $R(x)$.

Когда контур l недостаточно широк и близко подходит к отрезку $[a, b]$, могут существовать такие k и такие точки t на l , что отношение $\frac{x-x_k}{t-x_k}$ будет иметь модуль, больший единицы. Это может оказать влияние на величину $\frac{\omega(x)}{\omega(z)}$, и заключение о сходимости здесь получить значительно труднее. Чтобы она имела место, нужно, чтобы отношения $\frac{x-x_k}{t-x_k}$, меньшие единицы по модулю, оказывали бы большее влияние на $\frac{\omega(x)}{\omega(t)}$, чем отношения с модулем, большим единицы. Как учитывается влияние отношений обоих типов, будет показано в параграфе о сходимости интерполяционных процессов.

§ 4.4. НЕКОТОРЫЕ ПРАВИЛА ИНТЕРПОЛИРОВАНИЯ ПРИ РАВНООТСТОЯЩИХ ЗНАЧЕНИЯХ АРГУМЕНТА

В случае равноотстоящих узлов $x_0, x_1=x_0+h, \dots, x_n=x_0+nh, \dots$, который встречается в вычислительной практике особенно часто, правила интерполирования и соответствующие им формулы значительно упрощаются. Этот случай, ввиду его практической важности, привлекал к себе особенно много внимания, и здесь было построено чрезвычайно большое число правил интерполирования, предназначенных часто для узких целей. Мы ознакомимся только с теми из них, которые особенно часто применяются в вычислениях.

4.4.1. Правила для интерполирования в начале и конце таблицы

Предположим, что в равноотстоящих точках $x_k=x_0+kh$ ($k=0, 1, 2, \dots$) известны значения $f(x_k)=f(x_0+kh)=y_k$ функции $y(x)=f(x)$ и нам нужно интерполировать ее «вблизи» точки x_0 . Для интерполирования тогда естественно привлекать узлы x_k в том порядке, как они идут в таблице: $x_0, x_0+h, x_0+2h, \dots$

Применим для интерполирования правило Ньютона (4.3.6)

$$f(x) = f(x_0) + (x-x_0)f(x_0, x_0+h) + (x-x_0)(x-x_0-h)f(x_0, x_0+h, x_0+2h) + \dots + (x-x_0)(x-x_0-h) \dots (x-x_0-(k-1)h)f(x_0, x_0+h, \dots, x_0+kh) + R_k.$$

Примем прежде всего во внимание выражения (4.2.9) разностных отношений в равноотстоящих точках через конечные разности:

$$f(x_0) = y_0, f(x_0, x_0+h) = \frac{\Delta y_0}{1! h}, f(x_0, x_0+h, x_0+2h) = \frac{\Delta^2 y_0}{2! h^2}, \dots$$

Кроме того, введем новую переменную t , положив $x = x_0 + th$, $t = \frac{x - x_0}{h}$. Переменная t имеет смысл числа шагов h от x_0 до x .

$$\begin{aligned} x - x_0 &= th, (x - x_0)(x - x_0 - h) = t(t-1)h^2, \\ (x - x_0)(x - x_0 - h)(x - x_0 - 2h) &= t(t-1)(t-2)h^3, \dots \end{aligned}$$

Если внести все указанные величины в выражение для $f(x)$, получим правило Ньютона для интерполирования в начале таблицы

$$\begin{aligned} y(x_0 + th) &= y_0 + \frac{t}{1!} \Delta y_0 + \frac{t(t-1)}{2!} \Delta^2 y_0 + \frac{t(t-1)(t-2)}{3!} \Delta^3 y_0 + \\ &+ \dots + \frac{t(t-1)\dots(t-k+1)}{k!} \Delta^k y_0 + R_k. \end{aligned} \quad (4.4.1)$$

Считая $y = f(x)$ $k+1$ раз непрерывно дифференцируемой на отрезке, где лежат точки $x, x_0, x_0+h, \dots, x_0+kh$, применим для остатка R_k лагранжево представление (4.3.10).

$$\begin{aligned} \omega(x) &= (x - x_0)(x - x_0 - h)\dots(x - x_0 - kh) = h^{k+1}t(t-1)\dots(t-k), \\ R_k &= h^{k+1} \frac{t(t-1)\dots(t-k)}{(k+1)!} f^{(k+1)}(\xi), \end{aligned} \quad (4.4.2)$$

где ξ есть точка указанного выше отрезка. В частности, если точка интерполирования лежит между x_0 и x_0+kh , то ξ тоже принадлежит отрезку $[x_0, x_0+kh]$.

Предположим теперь, что мы находимся в конце таблицы с узлами $\dots, x_n - 3h, x_n - 2h, x_n - h, x_n$ и пусть точка интерполирования лежит вблизи x_n или где угодно справа от нее. При интерполировании в этом случае узлы следует брать в порядке их удаленности от x_n : $x_n, x_n - h, x_n - 2h, \dots$ Правило Ньютона (4.3.6) в этом случае запишется в таком виде:

$$\begin{aligned} f(x) &= f(x_n) + (x - x_n)f(x_n, x_n - h) + \\ &+ (x - x_n)(x - x_n + h)f(x_n, x_n - h, x_n - 2h) + \dots + \\ &+ (x - x_n)(x - x_n + h)\dots(x - x_n + (k-1)h)f(x_n, x_n - h, \dots, x_n - kh) + R_k. \end{aligned}$$

Если вновь принять во внимание, что

$$f(x_n) = y_n, \quad f(x_n, x_n - h) = f(x_n - h, x_n) = \frac{\Delta y_{n-1}}{1! h},$$

$$f(x_n, x_n - h, x_n - 2h) = f(x_n - 2h, x_n - h, x_n) = \frac{\Delta^2 y_{n-2}}{2! h^2}, \dots$$

и ввести переменную t , положив $x = x_n + th$, мы получим правило Ньютона для интерполирования в конце таблицы:

$$y(x_n + th) = y_n + \frac{t}{1!} \Delta y_{n-1} + \frac{t(t+1)}{2!} \Delta^2 y_{n-2} + \frac{t(t+1)(t+2)}{3!} \Delta^3 y_{n-3} +$$

$$+ \dots + \frac{t(t+1) \dots (t+k-1)}{k!} \Delta^k y_{n-k} + R_k, \quad (4.4.3)$$

остаток которого равен

$$R_k = h^{k+1} \frac{t(t+1) \dots (t+k)}{(k+1)!} f^{(k+1)}(\xi), \quad (4.4.4)$$

при этом точка ξ лежит на отрезке, содержащем $x_n - kh, \dots, x_n, x$.

4.4.2. Правила интерполирования внутри таблицы

Пусть x_n есть внутренний узел таблицы. Предположим, что точка интерполирования x лежит вблизи x_n с той или другой стороны. Табличные точки для интерполирования здесь разумно привлекать в следующем порядке: сначала взять x_n , затем брать пары точек $(x_n + h, x_n - h)$, $(x_n + 2h, x_n - 2h)$, \dots , $(x_n + kh, x_n - kh)$.

Число взятых узлов будет нечетным и равным $2k+1$.

Правило Ньютона (4.3.6) при таком порядке узлов запишется так:

$$f(x) = f(x_n) + (x - x_n) f(x_n, x_n + h) +$$

$$+ (x - x_n)(x - x_n - h) f(x_n, x_n + h, x_n - h) +$$

$$+ (x - x_n)(x - x_n - h)(x - x_n + h) f(x_n, x_n + h, x_n - h, x_n + 2h) + \dots +$$

$$+ (x - x_n)(x - x_n - h) \dots (x - x_n + kh - h) f(x_n, x_n + h, \dots, x_n + kh) +$$

$$+ (x - x_n)(x - x_n - h) \dots (x - x_n + kh - h)(x - x_n - kh) f(x_n, x_n + h, \dots,$$

$$x_n + kh, x_n - kh) + R_{2k},$$

$$R_{2k} = \frac{(x - x_n)(x - x_n - h)(x - x_n + h) \dots (x - x_n - kh)(x - x_n + kh)}{(2k+1)!} f^{(2k+1)}(\xi).$$

Вновь, если заменить x , введя переменную $t = \frac{x-x_n}{h}$, $x = x_n + th$, и подставить выражения разностных отношений через конечные разности

$$f(x_n) = y_n, \quad f(x_n, x_n + h) = \frac{\Delta y_n}{1! h},$$

$$f(x_n, x_n + h, x_n - h) = f(x_n - h, x_n, x_n + h) = \frac{\Delta^2 y_{n-1}}{2! h^2},$$

$$f(x_n, x_n + h, x_n - h, x_n + 2h) = f(x_n - h, x_n, x_n + h, x_n + 2h) = \frac{\Delta^3 y_{n-1}}{3! h^3}, \dots,$$

мы получим это равенство в виде:

$$\begin{aligned} y(x_n + th) = & y_n + \frac{t}{1!} \Delta y_n + \frac{t(t-1)}{2!} \Delta^2 y_{n-1} + \frac{t(t-1)(t+1)}{3!} \Delta^3 y_{n-1} + \\ & + \frac{t(t-1)(t+1)(t-2)}{4!} \Delta^4 y_{n-2} + \dots + \\ & + \frac{1}{(2k-1)!} (t+k-1) \dots t(t+1) \dots (t-k+1) \Delta^{2k-1} y_{n-k+1} + \\ & + \frac{1}{(2k)!} (t+k-1) \dots t \dots (t-k) \Delta^{2k} y_{n-k+1} + R_{2k}. \end{aligned}$$

Чтобы придать членам правой части симметричный вид, приведем сначала равенство к форме

$$\begin{aligned} y(x_n + th) = & y_n + t \left[\Delta y_n - \frac{1}{2} \Delta^2 y_{n-1} \right] + \frac{t^2}{2!} \Delta^2 y_{n-1} + \\ & + \frac{t(t^2-1^2)}{3!} \left[\Delta^3 y_{n-1} - \frac{1}{2} \Delta^4 y_{n-2} \right] + \dots + \\ & + \frac{t(t^2-1^2) \dots (t^2-(k-1)^2)}{(2k-1)!} \left[\Delta^{2k-1} y_{n-k+1} - \frac{1}{2} \Delta^{2k} y_{n-k} \right] + \\ & + \frac{t^2(t^2-1^2) \dots (t^2-(k-1)^2)}{(2k)!} \Delta^{2k} y_{n-k} + R_{2k}. \end{aligned}$$

Преобразуем, наконец, прямоугольные скобки, исключив в них разности четного порядка при помощи равенств

$$\Delta^2 y_{n-1} = \Delta y_n - \Delta y_{n-1}, \quad \Delta^4 y_{n-2} = \Delta^3 y_{n-1} - \Delta^3 y_{n-2}, \dots$$

В результате получим правило Ньютона — Стирлинга

$$\begin{aligned} y(x_n + th) = & y_n + \frac{t}{1!} \frac{\Delta y_n + \Delta y_{n-1}}{2} + \frac{t^2}{2!} \Delta^2 y_{n-1} + \\ & + \frac{t(t^2 - 1^2)}{3!} \frac{\Delta^3 y_{n-1} + \Delta^3 y_{n-2}}{2} + \frac{t^2(t^2 - 1^2)}{4!} \Delta^4 y_{n-2} + \dots + \\ & + \frac{t(t^2 - 1^2) \dots (t^2 - (k-1)^2)}{(2k-1)!} \frac{\Delta^{2k-1} y_{n-k+1} + \Delta^{2k-1} y_{n-k}}{2} + \\ & + \frac{t^2(t^2 - 1^2) \dots (t^2 - (k-1)^2)}{(2k)!} \Delta^{2k} y_{n-k} + R_{2k}, \end{aligned} \quad (4.4.5)$$

$$R_{2k} = h^{2k+1} \frac{t(t^2 - 1^2)(t^2 - 2^2) \dots (t^2 - k^2)}{(2k+1)^2} f^{(2k+1)}(\xi),$$

где ξ есть точка, принадлежащая отрезку, содержащему x_{n-k} , x_{n+k} и x .

Последнее правило, на котором мы остановимся, называется правилом Ньютона — Бесселя и предназначается для интерполирования в том случае, когда точка x лежит вблизи середины между табличными значениями. Пусть это будут значения x_n и x_{n+1} .

Соображения симметрии побуждают строить интерполяционное правило со следующим порядком привлечения узлов: сначала берется пара узлов $(x_n, x_n + h)$, затем пары $(x_n - h, x_n + 2h)$, $(x_n - 2h, x_n + 3h)$, \dots , $(x_n - kh + h, x_n + kh)$. Число узлов является четным. Правило Ньютона (4.3.6) при таком расположении узлов будет иметь следующий вид:

$$\begin{aligned} f(x) = & f(x_n) + (x - x_n) f(x_n, x_n + h) + \\ & + (x - x_n)(x - x_n - h) f(x_n, x_n + h, x_n - h) + \\ & + (x - x_n + h)(x - x_n)(x - x_n - h) f(x_n, x_n + h, x_n - h, x_n + 2h) + \dots + \\ & + (x - x_n + kh - 2h) \dots (x - x_n - kh + h) f(x_n, x_n + h, \dots, x_n + kh - h, \\ & x_n - kh + h) + (x - x_n + kh - h) \dots (x - x_n - kh + h) f(x_n, x_n + h, \dots, \\ & x_n - kh + h, x_n + kh) + R_{2k-1}, \\ R_{2k-1} = & \frac{(x - x_n + kh - h) \dots (x - x_n - kh)}{(2k)!} f^{(2k)}(\xi). \end{aligned}$$

После замены $x = x_n + th$, приняв во внимание равенства

$$f(x_n) = y_n, \quad f(x_n, x_n + h) = \frac{\Delta y_n}{1! h}, \quad f(x_n, x_n + h, x_n - h) = \frac{\Delta^2 y_{n-1}}{2! h^2},$$

$$f(x_n, x_n + h, x_n - h, x_n + 2h) = \frac{\Delta^3 y_{n-1}}{3! h^3}, \dots,$$

найдем

$$\begin{aligned} y(x_n + th) = & y_n + \frac{t}{1!} \Delta y_n + \frac{t(t-1)}{2!} \Delta^2 y_{n-1} + \frac{(t+1)t(t-1)}{3!} \Delta^3 y_{n-1} + \\ & + \frac{(t+1)t(t-1)(t-2)}{4!} \Delta^4 y_{n-2} + \dots + \frac{(t+k-2) \dots (t-k+1)}{(2k-2)!} \Delta^{2k-2} y_{n-k+1} + \\ & + \frac{(t+k-1) \dots (t-k+1)}{(2k-1)!} \Delta^{2k-1} y_{n-k+1} + R_{2k-1}. \end{aligned}$$

Для приведения членов правой части к виду, симметричному относительно точки $x_n + \frac{1}{2}h$, отделим от четных разностей половины их значений

$$\frac{1}{2} y_n, \quad \frac{1}{2} \Delta^2 y_{n-1}, \quad \frac{1}{2} \Delta^4 y_{n-2}, \dots$$

и заменим эти значения при помощи тождеств

$$\frac{1}{2} y_n = \frac{1}{2} (y_{n+1} - \Delta y_n), \quad \frac{1}{2} \Delta^2 y_{n-1} = \frac{1}{2} (\Delta^2 y_n - \Delta^3 y_{n-1}),$$

$$\frac{1}{2} \Delta^4 y_{n-2} = \frac{1}{2} (\Delta^4 y_{n-1} - \Delta^5 y_{n-2}), \dots$$

После приведения подобных членов получим интерполяционное правило Ньютона — Бесселя:

$$\begin{aligned} y(x_n + th) = & \frac{y_n + y_{n+1}}{2} + \frac{t - \frac{1}{2}}{1!} \Delta y_n + \frac{t(t-1)}{2!} \frac{\Delta^2 y_{n-1} + \Delta^2 y_n}{2} + \\ & + \frac{\left(t - \frac{1}{2}\right) t(t-1)}{3!} \Delta^3 y_{n-1} + \frac{(t+1)t(t-1)(t-2)}{4!} \frac{\Delta^4 y_{n-2} + \Delta^4 y_{n-1}}{2} + \end{aligned}$$

$$\begin{aligned}
& + \dots + \frac{(t+k-2) \dots (t-k+1)}{(2k-2)!} \frac{\Delta^{2k-2} y_{n-k+1} + \Delta^{2k-2} y_{n-k+2}}{2} + \\
& + \frac{\left(t - \frac{1}{2}\right) (t+k-2) \dots (t-k+1)}{(2k-1)!} \Delta^{2k-1} y_{n-k+1} + R_{2k-1}, \quad (4.4.6)
\end{aligned}$$

$$R_{2k-1} = h^{2k} \frac{(t+k-1) \dots (t-k)}{(2k)!} f^{(2k)}(\xi),$$

где ξ есть некоторая точка отрезка, содержащего $x_n - kh + h$, $x_n + kh$, x .

§ 4.5. ПРИЛОЖЕНИЕ ИНТЕРПОЛИРОВАНИЯ К ЧИСЛЕННОМУ НАХОЖДЕНИЮ ПРОИЗВОДНЫХ

4.5.1. Об интерполяционном правиле вычисления производной от функции, заданной таблично

Такая задача может быть решена только приближенно. Наряду с установлением правил вычисления производных одной из основных задач здесь является оценка погрешности, которая допускается при вычислении.

Предположим, что для $(n+1)$ -кратно непрерывно дифференцируемой функции f в различных точках x_0, x_1, \dots, x_n отрезка $[a, b]$ известны ее значения $y_0 = f(x_0), y_1 = f(x_1), \dots, y_n = f(x_n)$. По этим исходным данным нужно найти значение производной порядка m от $f(x)$ в любой точке $x \in [a, b]$. Одно из возможных решений задачи состоит в следующем. По исходным данным выполняют алгебраическое интерполирование f . Пусть $P(x)$ есть интерполирующий многочлен, имеющий степень не выше n , и $R(x)$ — остаток интерполирования, так что

$$f(x) = P(x) + R(x).$$

Вычислим производную порядка m от обоих членов равенства:

$$f^{(m)}(x) = P^{(m)}(x) + R^{(m)}(x). \quad (4.5.1)$$

Если пренебречь здесь величиной $R^{(m)}(x)$, получим численное правило для нахождения нужной производной:

$$f^{(m)}(x) \approx P^{(m)}(x), \quad (4.5.2)$$

погрешность которого равна $R^{(m)}(x)$. При применении правила мы должны, очевидно, считать $m \leq n$, так как все производные от P порядка, большего n , равны нулю тождественно.

Вычисление производной $P^{(m)}(x)$ от многочлена P принципиальных трудностей не представляет. Наибольший интерес здесь имеет лишь важный в практике технический вопрос о приведении вычислений к виду, возможно более удобному в работе. Отложим эту задачу на некоторое время и займемся сейчас остатком $R^{(m)}(x)$.

$R(x)$ есть $n+1$ раз непрерывно дифференцируемая функция на $[a, b]$, обращающаяся в нуль в $n+1$ узлах x_0, x_1, \dots, x_n . Чтобы сделать более наглядными рассуждения, будем считать, что узлы x_k перенумерованы в порядке роста координат: $x_k < x_{k+1}$ ($k=0, 1, \dots, n-1$).

Рассмотрим вспомогательную функцию

$$\varphi(t) = R(t) - k \frac{\omega(t)}{(n+1)!}, \quad \omega(t) = (t-x_0)(t-x_1)\dots(t-x_n), \quad (4.5.3)$$

где k есть постоянная, которую мы выберем позже. Аргумент функции φ мы обозначили буквой t , чтобы отличать его от точки x , в которой вычисляется $f^{(m)}(x)$. Функция $\varphi(t)$ обращается в нуль в узлах x_k и имеет на отрезке $[x_0, x_n]$ по меньшей мере $n+1$ различных нулей. Между каждой парой узлов (x_k, x_{k+1}) φ' будет иметь по меньшей мере один нуль, и число нулей φ' внутри (x_0, x_n) будет не меньше n и т. д. Производная порядка m

$$\varphi^{(m)}(t) = R^{(m)}(t) - \frac{k}{(n+1)!} \omega^{(m)}(t)$$

будет иметь внутри (x_0, x_n) по меньшей мере $n+1-m$ различных нулей.

Займемся выбором числа k . Предположим, что точка интерполирования x не лежит внутри отрезка $[x_1, x_n]$, а располагается вне его или на одном из его концов. Потребуем, чтобы в точке $t=x$ выполнялось равенство

$$\varphi^{(m)}(x) = R^{(m)}(x) - \frac{k}{(n+1)!} \omega^{(m)}(x) = 0 \quad (m \geq 1).$$

Так как все нули $\omega^{(m)}(x)$ лежат внутри $[x_0, x_n]$, то $\omega^{(m)}(x) \neq 0$ и число k может быть найдено.* Кроме того, так как точка x по предположению не лежит внутри $[x_0, x_n]$, то $\varphi^{(m)}(t)$ имеет на отрезке, содержащем x , x_0, \dots, x_n , по меньшей мере $n+2-m$ разных нулей. Производная $\varphi^{(m+1)}(t)$ будет иметь внутри этого отрезка по меньшей мере $n+1-m$ различных нулей и т. д. и производная $\varphi^{(n+1)}(t)$ должна иметь внутри отрезка,

*) $\omega'(x)$ есть многочлен степени n , имеющий внутри каждого из отрезков (x_k, x_{k+1}) ($k=0, 1, \dots, n-1$) по одному нулю и не имеющий никаких других нулей. $\omega''(x)$ есть многочлен степени $n-1$, имеющий по одному нулю между каждыми соседними нулями $\omega'(x)$. Таких нулей $n-1$, и все они лежат внутри $[x_0, x_n]$. Никаких других нулей у $\omega''(x)$ нет и т. д.

не меньше одного нуля. Назовем его ξ . Ввиду $\omega^{(n+1)}(t) = (n+1)!$ и $R^{(n+1)} = f^{(n+1)} - P^{(n+1)} = f^{(n+1)}$, при $t = \xi$ должно быть

$$\varphi^{(n+1)}(\xi) = R^{(n+1)}(\xi) - k = f^{(n+1)}(\xi) - k = 0, \quad k = f^{(n+1)}(\xi).$$

Из уравнения для определения k следует

$$f^{(m)}(x) - P^{(m)}(x) = R^{(m)}(x) = \frac{1}{(n+1)!} \omega^{(m)}(x) f^{(n+1)}(\xi), \quad (4.5.4)$$

что позволяет высказать приводимую ниже теорему.

Теорема 1. Пусть на отрезке $[a, b]$, содержащем x и узлы x_0, \dots, x_n , функция f имеет непрерывную производную и x не лежит между x_0 и x_n . Тогда на указанном выше отрезке существует такая точка ξ , что для погрешности $R^{(m)}(x)$ вычисления производной $f^{(m)}(x)$ верно равенство (4.5.4).

Если точка x лежит внутри отрезка $[x_0, x_1]$, то наши рассуждения становятся неправомерными в двух пунктах. Может оказаться, что в точке x производная $\omega^{(m)}(x) = 0$ и уравнение для k станет неразрешимым или неопределенным. Если уравнение для k разрешимо, то нельзя поручиться за то, что корень $t = x$ для $\varphi^{(m)}(t)$ будет отличным от тех $n+1-m$ корней, существование которых было доказано при помощи теоремы Ролля, и может оказаться, что $\varphi^{(m)}(t)$ имеет меньше чем $n+2-m$ корней. Поэтому при расположении точки x между x_0 и x_n мы должны считаться с возможностью исключительных случаев, когда исследуемая погрешность $R^{(m)}(x)$ может не иметь представления вида (4.5.4). Какое можно построить представление $R^{(m)}(x)$, верное при всяком расположении точки x , мы увидим несколькими строками ниже. Сейчас же рассмотрим одну задачу вычисления производной первого порядка f' при расположении точки x внутри $[x_0, x_n]$, в которой остается верным представление погрешности (4.5.4).

Пусть x_i есть любой из узлов и мы хотим вычислить $f'(x_i)$. Приближенно положим $f'(x_i) \approx P'(x_i)$ и рассмотрим погрешность $R'(x_i) = f'(x_i) - P'(x_i)$. Вспомогательная функция $\varphi(t)$, как отмечалось выше, обращается в нуль в каждом из узлов x_k ($k=0, 1, \dots, n$), и $\varphi'(t)$ будет иметь не менее одного нуля внутри каждого из отрезков $[x_k, x_{k+1}]$ ($k=0, 1, \dots, n$). Таких нулей будет по меньшей мере n штук.

Выберем теперь k так, чтобы точка $t = x_i$ также была нулем φ' :

$$\varphi'(x_i) = R'(x_i) - \frac{k}{(n+1)!} \omega'(x_i) = 0,$$

$\omega'(x_i) = \prod_{j \neq i} (x_i - x_j) \neq 0$ и число k может быть найдено. Нуль $t = x_i$ отли-

чен от нулей, лежащих внутри отрезков $[x_k, x_{k+1}]$, и, таким образом, φ' будет иметь не меньше $n+1$ различных нулей на $[x_0, x_n]$. Отсюда следует, что $\varphi^{(n+1)}(t)$ будет иметь на $[x_0, x_n]$ по крайней мере один нуль:

$$\varphi^{(n+1)}(\xi) = R^{(n+1)}(\xi) - k = 0, \quad k = R^{(n+1)}(\xi) = f^{(n+1)}(\xi).$$

Далее остается только повторить рассуждения, проделанные выше для производной любого порядка m , и мы придем к заключению, что погрешность $R'(x_i) = f'(x_i) - P'(x_i)$ приближенного вычисления производной в узле x_i по правилу $f'(x_i) \approx P'(x_i)$ имеет представление

$$R'(x_i) = \frac{\omega'(x_i)}{(n+1)!} f^{(n+1)}(\xi) \quad (x_0 < \xi < x_n). \quad (4.5.5)$$

Теперь укажем представление погрешности $R^{(m)}(x)$, верное при всяком положении точки x на $[a, b]$. Как и выше, будем считать f имеющей непрерывную на $[a, b]$ производную $f^{(n+1)}$. По теореме Тейлора для f верно равенство

$$f(x) = c_0 + c_1(x-a) + \dots + c_n(x-a)^n + \int_a^x \varphi(t) \frac{(x-t)^n}{n!} dt, \quad (4.5.6)$$

$$c_i = \frac{1}{i!} f^{(i)}(a) \quad (i=0, 1, \dots, n), \quad \varphi(t) = f^{(n+1)}(t), \quad a \leq x \leq b,$$

которое, по существу дела, является структурной формулой или параметрическим представлением множества функций f , $n+1$ раз непрерывно дифференцируемых на $[a, b]$, так как всякая функция f этого множества представима этой формулой и, наоборот, при любых значениях c_i ($i=0, 1, \dots, n$) и всякой функции φ , непрерывной на $[a, b]$, функция f , определенная равенством (4.5.6), является непрерывно дифференцируемой $n+1$ раз на $[a, b]$.

В (4.5.6) для наших целей удобнее заменить интеграл с переменной верхней границей на интеграл по отрезку $[a, b]$, что можно сделать, если ввести «гасящую» функцию, позволяющую уничтожить лишние участки интегрирования. Определим $E(x)$ равенством

$$E(x) = \begin{cases} 1 & \text{при } x > 0, \\ \frac{1}{2} & \text{при } x = 0, \\ 0 & \text{при } x < 0. \end{cases} \quad (4.5.6')$$

Как сразу же видно, (4.5.6) может быть записано в форме

$$f(x) = \sum_{i=0}^n c_i (x-a)^i + \int_a^b f^{(n+1)}(t) E(x-t) \frac{(x-t)^n}{n!} dt. \quad (4.5.7)$$

Остаток

$$R(f; x) = f(x) - \sum_{k=0}^n \frac{\omega(x)}{(x-x_k)\omega'(x_k)} f(x_k)$$

интерполирования функции f будет равен остатку интерполирования интегрального члена в (4.5.7), так как многочлен $\sum_{i=0}^n c_i (x-a)^i$, имеющий степень не выше n , интерполируется точно:

$$\begin{aligned} R(f; x) &= \int_a^b f^{(n+1)}(t) \left\{ E(x-t) \frac{(x-t)^n}{n!} - \right. \\ &- \sum_{k=0}^n \frac{\omega(x)}{(x-x_k)\omega'(x_k)} E(x_k-t) \frac{(x_k-t)^n}{n!} \left. \right\} dt = \\ &= \int_a^b f^{(n+1)}(t) R \left[E(x-t) \frac{(x-t)^n}{n!}; x \right] dt. \end{aligned}$$

При вычислении производной порядка $m \leq n$ по x последний интеграл можно дифференцировать по x под его знаком и для интересующей нас погрешности получится равенство

$$R^{(m)}(f; x) = \int_a^b f^{(n+1)}(t) R_{x^m}^{(m)} \left[E(x-t) \frac{(x-t)^n}{n!}; x \right] dt. \quad (4.5.8)$$

Это выражение для $R^{(m)}(f; x)$ остается, очевидно, верным при всяком расположении точки x на $[a, b]$. Равенство же (4.5.4) получится отсюда, если применить к интегралу теорему о среднем значении, что возможно, вообще говоря, лишь в том случае, когда ядро интеграла

$$R_{x^m}^{(m)} \left[E(x-t) \frac{(x-t)^n}{n!}; x \right]$$

сохраняет, как функция от t , знак на отрезке $a \leq t \leq b$.

Если рассмотреть множество функций f , имеющих на $[a, b]$ непрерывную производную порядка $n+1$, удовлетворяющую условию $|f^{(n+1)}(x)| \leq M_{n+1}$, то для него остаток $R_{x^n}^{(m)}(f; x)$ имеет следующую точную оценку:

$$|R_{x^n}^{(m)}(f; x)| \leq M_{n+1} \int_a^b \left| R_{x^n}^{(m)} \left[E(x-t) \frac{(x-t)^n}{n!}; x \right] \right| dt. \quad (4.5.9)$$

4.5.2. Некоторые частные правила вычисления производных

При вычислениях на быстродействующих вычислительных машинах, когда мы заинтересованы скорее не в малом количестве арифметических операций, а в экономии элементов памяти и в простоте программирования, достаточно выгодным является использование представления интерполирующего многочлена $P(x)$ в форме Лагранжа (4.3.5). Например, при вычислении первой производной, когда x не совпадает ни с одним узлом x_k ($k=0, 1, \dots, n$), можно воспользоваться следующим выражением для P' :

$$\begin{aligned} f'(x) &\approx P'(x) = \omega'(x) \sum_{k=0}^n \frac{1}{(x-x_k)\omega'(x_k)} f(x_k) - \\ &- \omega(x) \sum_{k=0}^n \frac{1}{(x-x_k)^2\omega'(x_k)} f(x_k) = \\ &= \omega(x) \left[\sum_{i=0}^n \frac{1}{x-x_i} \sum_{k=0}^n \frac{1}{(x-x_k)\omega'(x_k)} f(x_k) - \right. \\ &\quad \left. - \sum_{k=0}^n \frac{1}{(x-x_k)^2\omega'(x_k)} f(x_k) \right]. \end{aligned} \quad (4.5.10)$$

Если же вычислению подлежит значение f' в узле x_j , то можно воспользоваться более простым равенством

$$f'(x_j) \approx P'(x_j) = \omega'(x_j) \sum_{k \neq j} \frac{1}{(x_j-x_k)\omega'(x_k)} f(x_k) + \frac{1}{2} \frac{\omega''(x_j)}{\omega'(x_j)} f(x_j). \quad (4.5.11)$$

При счете на настольных машинах без программного управления более рационально, по-видимому, воспользоваться ньютоновой формой

(4.3.6) многочлена $P(x)$. Обозначив для сокращения записи $x - x_k = \alpha_k$, можно придать $P(x)$ форму

$$P(x) = f(x_0) + \alpha_0 f(x_0, x_1) + \alpha_0 \alpha_1 f(x_0, x_1, x_2) + \alpha_0 \alpha_1 \alpha_2 f(x_0, x_1, x_2, x_3) + \\ + \dots + \alpha_0 \alpha_1 \dots \alpha_{n-1} f(x_0, x_1, \dots, x_n).$$

Вычисляя производные от обеих частей, получим приводимые ниже выражения для f', f'', \dots :

$$f'(x) \approx P'(x) = f(x_0, x_1) + (\alpha_0 + \alpha_1) f(x_0, x_1, x_2) + \\ + (\alpha_0 \alpha_1 + \alpha_1 \alpha_2 + \alpha_2 \alpha_0) f(x_0, x_1, x_2, x_3) + \dots, \\ \frac{1}{2!} f''(x) \approx \frac{1}{2!} P''(x) = f(x_0, x_1, x_2) + \\ + (\alpha_0 + \alpha_1 + \alpha_2) f(x_0, x_1, x_2, x_3) + \dots, \\ \frac{1}{3!} f'''(x) \approx \frac{1}{3!} P'''(x) = f(x_0, x_1, x_2, x_3) + \\ + (\alpha_0 + \alpha_1 + \alpha_2 + \alpha_3) f(x_0, x_1, x_2, x_3, x_4) + \dots, \\ \frac{1}{4!} f^{IV}(x) \approx \frac{1}{4!} P^{IV}(x) = f(x_0, x_1, x_2, x_3, x_4) + \\ + (\alpha_0 + \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4) f(x_0, x_1, x_2, x_3, x_4, x_5) + \dots$$

В частном случае, когда вычисляются значения производных f', f'', \dots в узле x_0 , в предыдущих равенствах нужно положить $x = x_0$, $\alpha_0 = 0$, $\alpha_i = x_0 - x_i$ ($i = 1, 2, \dots$). После этого получатся равенства:

$$f'(x_0) \approx f(x_0, x_1) + (x_0 - x_1) f(x_0, x_1, x_2) + \\ + (x_0 - x_1)(x_0 - x_2) f(x_0, x_1, x_2, x_3) + \dots, \\ \frac{1}{2!} f''(x_0) \approx f(x_0, x_1, x_2) + (2x_0 - x_1 - x_2) f(x_0, x_1, x_2, x_3) + \\ + [(x_0 - x_1)(x_0 - x_2) + (x_0 - x_1)(x_0 - x_3) + \\ + (x_0 - x_2)(x_0 - x_3)] f(x_0, x_1, x_2, x_3, x_4) + \dots, \\ \frac{1}{3!} f'''(x_0) \approx f(x_0, x_1, x_2, x_3) + (3x_0 - x_1 - x_2 - x_3) f(x_0, x_1, x_2, x_3, x_4) + \dots, \\ \frac{1}{4!} f^{IV}(x_0) \approx f(x_0, x_1, x_2, x_3, x_4) + \\ + (4x_0 - x_1 - x_2 - x_3 - x_4) f(x_0, x_1, x_2, x_3, x_4, x_5) + \dots, \\ \dots \dots \dots$$

Совершенно аналогично могут быть получены выражения производных через конечные разности в случае равноотстоящих узлов.*) Если, например, исходить из правила Ньютона для интерполирования в начале таблицы (4.4.1), получатся следующие выражения для производных:

$$hy'(x_0+th) = \Delta y_0 + \frac{2t-1}{2!} \Delta^2 y_0 + \frac{3t^2-6t+2}{3!} \Delta^3 y_0 + \\ + \frac{4t^3-18t^2+22t-6}{4!} \Delta^4 y_0 + \dots,$$

$$h^2 y''(x_0+th) = \Delta^2 y_0 + (t-1) \Delta^3 y_0 + \frac{6t^2-18t+11}{12} \Delta^4 y_0 + \dots,$$

$$h^3 y'''(x_0+th) = \Delta^3 y_0 + \frac{2t-3}{2} \Delta^4 y_0 + \dots,$$

$$\dots \dots \dots$$

При $x=x_0$ и $t=0$ будет

$$hy'(x_0) = \Delta y_0 - \frac{1}{2} \Delta^2 y_0 + \frac{1}{3} \Delta^3 y_0 - \frac{1}{4} \Delta^4 y_0 + \frac{1}{5} \Delta^5 y_0 - \frac{1}{6} \Delta^6 y_0 + \dots,$$

$$h^2 y''(x_0) = \Delta^2 y_0 - \Delta^3 y_0 + \frac{11}{12} \Delta^4 y_0 - \frac{5}{6} \Delta^5 y_0 + \frac{137}{180} \Delta^6 y_0 - \dots,$$

$$h^3 y'''(x_0) = \Delta^3 y_0 - \frac{3}{2} \Delta^4 y_0 + \frac{7}{4} \Delta^5 y_0 - \frac{15}{8} \Delta^6 y_0 + \dots,$$

$$h^4 y^{IV}(x_0) = \Delta^4 y_0 - 2\Delta^5 y_0 + \frac{17}{6} \Delta^6 y_0 - \dots,$$

$$h^5 y^V(x_0) = \Delta^5 y_0 - \frac{5}{2} \Delta^6 y_0 + \dots,$$

$$\dots \dots \dots$$

*) Достаточно большое число правил вычисления производных приведено в книгах; [5] и И. С. Березин, Н. П. Жидков. Методы вычислений, т. 1, М., 1966.

§ 4.6. ИНТЕРПОЛЯЦИОННЫЕ МЕТОДЫ РЕШЕНИЯ ЧИСЛЕННЫХ УРАВНЕНИЙ

4.6.1. Введение. Связь с задачей обратного интерполирования

В гл. 1, когда обсуждался вопрос об уточнении метода Ньютона, мы обращали внимание на то, что одним из возможных способов уточнения является интерполирование. Ознакомимся сейчас с одной из простейших форм интерполяционного метода и ограничимся изложением лишь наглядной стороны вопроса.

Сначала обратим внимание на связь интерполяционного метода с так называемой «задачей обратного интерполирования».

Пусть рассматривается некоторая функция $y=f(x)$, для которой известна таблица ее значений

$$\left. \begin{array}{cccccc} x_0 & x_1 & x_2 & \dots & x_{n-1} & x_n \\ y_0 & y_1 & y_2 & \dots & y_{n-1} & y_n \end{array} \right\} (y_k=f(x_k)). \quad (4.6.1)$$

В обычной задаче интерполирования рассматривается такой вопрос: дано табличное значение аргумента x и нужно найти соответствующее ему значение $y=f(x)$. Рассмотрим обратную задачу: пусть задано табличное значение y^* функции $y=f(x)$ и нужно найти, какому значению аргумента x оно отвечает. По существу дела, здесь мы ставим задачу о решении уравнения $f(x)=y^*$, в котором число y^* считается заданным, а аргумент x является неизвестной величиной. При этом функция $f(x)$ считается заданной не аналитически, а таблично и при разыскании x мы имеем право пользоваться только числами x_k , y_k , входящими в таблицу (4.6.1), и значением y^* .

К решению поставленной задачи можно идти двумя путями. Первый из них, от которого произошло название задачи, является особенно простым, но относится к функциям $f(x)$ частного типа, когда функциональная зависимость $y=f(x)$ является однозначно обратимой. Известно, что это имеет место, наверное, в том случае, когда $f(x)$ является монотонной (возрастающей или убывающей) функцией. $f(x)$ считается заданной таблично, и по таблице значений (4.6.1) легко проследить, будет ли действительно f обладать свойством монотонности. Допустим, что монотонность имеет место, и рассмотрим обратную функцию $x=F(y)$. Она задана той же таблицей значений (4.6.1), что и $f(x)$, с тем различием, что y_k — теперь значения аргумента, а x_k — соответствующие им значения функции.

Интерполируем $F(y)$ при помощи многочлена $\Pi(y)$ степени n : $F(y) \approx \Pi(y)$. Положив здесь $y=y^*$, найдем приближенно нужное нам значение x : $x \approx \Pi(y^*)$.

Второй путь нахождения x имеет более общее значение и применим ко всякой функции f . Но он более труден по вычислениям, так как требует решения алгебраического уравнения. Возвратимся к заданной функ-

ции $f(x) = y$ и проинтерполируем ее при помощи многочлена $P(x)$ степени n . Уравнение $f(x) = y^*$ заменим новым уравнением $P(x) \approx y^*$.

Обычно бывает, что значения x_k ($k=0, 1, \dots, n$) принадлежат малому участку и f изменяется на нем достаточно гладко. Кроме того, y^* является, как правило, близким к y_k ($k=0, 1, \dots, n$), и поэтому уравнение $P(x) \approx y^*$ будет мало отличаться от $f(x) = y^*$ и можно ожидать, что решение нового, приближенного уравнения будет близким к искомому значению x^* .)

Все изложенные соображения являются, разумеется, очень нестрогими, но они достаточны для наших целей и позволят показать, какое полезное значение имеет проблема обратного интерполирования в задаче решения численных уравнений. Возвратимся к этой последней задаче. Предположим, что дано уравнение $f(x) = 0$. Функция f задается аналитически, и мы имеем возможность вычислить значение f в любой точке вблизи решения.

Пусть каким-либо путем найдено несколько приближений к решению и составлена таблица (4.6.1). Нам нужно указать правило для нахождения следующего приближения x_{n+1} . Условия этой задачи весьма сходны с теми, в которых мы находились в проблеме обратного интерполирования. Нам дано уравнение $f(x) = 0$ и заданное заранее значение функции y^* здесь равно нулю. В качестве следующего приближения x_{n+1} естественно поэтому взять то значение x , которое получится по методу обратного интерполирования. В соответствии с этим перед нами открывается возможность построить x_{n+1} двумя методами, соответственно двум возможным путям обратного интерполирования.

Ниже мы проследим несколько более подробно каждый из двух методов, а сейчас выясним причины, по которым избранная нами форма интерполяционного метода была названа простейшей. Здесь дело в том, что для нахождения x_{n+1} мы выполняем интерполирование только по значениям функции и не пользуемся значениями производных. Напомним, что в методе Ньютона следующее приближение x_{n+1} находится по правилу

$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$ и при вычислениях ведется таблица вида

$$\left. \begin{array}{cccccc} x_0 & x_1 & & & x_{n-1} & x_n \\ y_0 & y_1 & \dots & & y_{n-1} & y_n \\ y'_0 & y'_1 & \dots & & y'_{n-1} & y'_n \end{array} \right\}$$

Если строить интерполяционные методы, уточняющие правило Ньютона, мы могли бы воспользоваться интерполированием по значениям

) Уравнение $P(x) \approx y^$ может иметь несколько решений, и из всех его решений мы должны выбрать то, которое лежит наиболее близко к x_i ($i=0, 1, \dots, n$). Такое решение чаще всего бывает единственным.

функции $y=f(x)$ и первой производной $y'=f'(x)$. Принципиально говоря, можно было бы строить интерполяционные методы решения уравнений с узлами любых кратностей. Мы же для облегчения изложения избрали самую простую форму метода, когда интерполирование выполняется только по значениям f , но хотим обратить внимание читателя на возможность его обобщений и уточнений.

4.6.2. Метод приближений, основанный на интерполировании обратной функции

Точное решение уравнения $f(x)=0$ обозначим x^* и предположим, что на некотором отрезке, содержащем x^* и все последовательные приближения x_k ($k=0, 1, \dots$), о которых будет говориться ниже, $f(x)$ имеет однозначную обратную функцию $F(x)$.

Пусть вычисления начаты, доведены до приближения x_n и составлена таблица (4.6.1). Рассмотрим обратную функцию $x=F(y)$ и проинтерполируем ее по $k+1$ ($k \leq n$) значениям, которые она принимает в узлах $y_n, y_{n-1}, \dots, y_{n-k}$. Интерполирующий многочлен запишем в форме Лагранжа, позволяющей представить все величины в легко обозримом виде, хотя эта форма и не является, по-видимому, самой удобной для вычислений,*)

$$F(y) \approx \Pi(y) = \sum_{j=0}^k \frac{\Omega(y)}{(y-y_{n-j})\Omega'(y_{n-j})} x_{n-j}, \quad (4.6.2)$$

$$\Omega(y) = (y-y_n)(y-y_{n-1}) \dots (y-y_{n-k}).$$

Положим здесь $y=0$ и значение $\Pi(0)$ примем за x_{n+1} :

$$x_{n+1} = \Pi(0) = \sum_{j=0}^k \frac{\Omega(0)}{-y_{n-j}\Omega'(y_{n-j})} x_{n-j}. \quad (4.6.3)$$

Затем вычисляем $f(x_{n+1})=y_{n+1}$, добавляем еще один столбец в таблицу (4.6.1) и переходим к нахождению x_{n+2} .

Рассмотрим теперь погрешность приближения к решению $\varepsilon_n = x^* - x_n$ и постараемся выяснить наглядную картину изменения ε_n , когда x_n близко к x^* . Будем считать, что обратная функция F непрерывно дифференцируема $k+1$ раз. Это, наверное, будет так, если $f(x)$ непрерывно дифференцируема $k+1$ раз и $f'(x) \neq 0$.

*) При вычислениях на машинах без программного управления более удобным, вероятно, было бы воспользоваться интерполяционной формулой Ньютона, рассчитанной на интерполирование в конце таблицы (4.6.1): $\Pi(y) = F(y_n) + (y-y_n)F'(y_n) + \frac{(y-y_n)^2}{2!}F''(y_n) + \dots + \frac{(y-y_n)^{k+1}}{(k+1)!}F^{(k+1)}(y_n) + \dots$

Остаток интерполирования $F(y)$ равен

$$r(y) = F(y) - \Pi(y) = \frac{\Omega(y)}{(k+1)!} F^{(k+1)}(\eta),$$

где η есть промежуточная точка на отрезке, содержащем y_{n-i} ($i=0, 1, \dots, k$) и y .

Положим здесь $y=0$ и примем во внимание, что $F(0)=x^*$, $\Pi(0)=x_{n+1}$:

$$\begin{aligned} \varepsilon_{n+1} &= x^* - x_{n+1} = F(0) - \Pi(0) = r(0) = \frac{\Omega(0)}{(k+1)!} F^{(k+1)}(\eta_0) = \\ &= \frac{(-1)^{k+1}}{(k+1)!} \prod_{i=0}^k y_{n-i} F^{(k+1)}(\eta_0) = \frac{(-1)^{k+1}}{(k+1)!} \prod_{i=0}^k f(x_{n-i}) F^{(k+1)}(\eta_0). \end{aligned}$$

Под η_0 здесь понимается точка отрезка, содержащего y_{n-i} ($i=0, 1, \dots, k$) и точку нуль.

Ввиду $f(x^*)=0$ будет

$$f(x_m) = f(x_m) - f(x^*) = -\varepsilon_m f'(x^* - \theta_m \varepsilon_m) \quad (0 < \theta_m < 1),$$

$$\varepsilon_{n+1} = \frac{1}{(k+1)!} \prod_{i=0}^k \varepsilon_{n-i} \prod_{i=0}^k f'(x^* - \theta_{n-i} \varepsilon_{n-i}) F^{(k+1)}(\eta_0). \quad (4.6.4)$$

Чтобы сделать наглядным закон изменения погрешности ε_n , сохраним в правой части лишь главный член, производя замену $f'(x^* - \theta_{n-i} \varepsilon_{n-i}) \approx f'(x^*)$ и $F^{(k+1)}(\eta_0) \approx F^{(k+1)}(0)$. После этого получим приближенное равенство, дающее достаточно простую картину изменения ε_n вблизи решения x^* :

$$\varepsilon_{n+1} \approx \frac{1}{(k+1)!} [f'(x^*)]^{k+1} F^{(k+1)}(0) \varepsilon_n \varepsilon_{n-1} \dots \varepsilon_{n-k}. \quad (4.6.5)$$

Равенство (4.6.5) позволяет думать, что если f вблизи решения x^* имеет непрерывную производную порядка $k+1$ и первую производную, отличную от нуля, а приближения x_0, x_1, \dots, x_k взяты достаточно близкими к точному решению, то вычислительный процесс, определяемый правилом (4.6.3), будет сходиться к x^* весьма быстро, со скоростью, указанной в (4.6.5).

4.6.3. Замена точного уравнения $f(x)=0$ приближенным, полученным интерполированием f

Метод нахождения x^* , о котором здесь будет идти речь, основан на замене заданного уравнения более, вообще говоря, простым алгебраическим уравнением. Рассмотрим $k+1$ значений, принимаемых функцией f в узлах $x_n, x_{n-1}, \dots, x_{n-k}$, и проинтерполируем ее по этим значениям при помощи многочлена $P(x)$ степени $\leq n$.

$$f(x) = P(x) + R(x), \quad R(x) = \frac{\omega(x)}{(k+1)!} f^{(k+1)}(\xi),$$

$$\omega(x) = (x-x_n)(x-x_{n-1}) \dots (x-x_{n-k}),$$

ξ — принадлежит отрезку, содержащему точки x_n, \dots, x_{n-k} , x . Отбрасывая остаток интерполирования $R(x)$, заменим заданное уравнение $f(x)=0$ «близким» к нему вспомогательным алгебраическим уравнением $P(x)=0$. Новое уравнение может иметь несколько решений, и из них выбирают то, которое будет ближайшим к месту расположения узлов x_n, \dots, x_{n-k} , или, если этот признак окажется недостаточно определенным, ближайшим к последнему известному приближению x_n .

Решать уравнение $P(x)=0$ удобно методом итерации в следующем виде. Воспользуемся формулой Ньютона для интерполирования в конце таблицы и запишем уравнение в виде

$$P(x) = f(x_n) + (x-x_n)f(x_n, x_{n-1}) + \\ + (x-x_n)(x-x_{n-1})f(x_n, x_{n-1}, x_{n-2}) + \dots = 0.$$

Перенесем член $(x-x_n)f(x_n, x_{n-1})$ в другую часть равенства и разделим обе части на $-f(x_n, x_{n-1})$:

$$x = x_n - \frac{f(x_n)}{f(x_n, x_{n-1})} - (x-x_n)(x-x_{n-1}) \frac{f(x_n, x_{n-1}, x_{n-2})}{f(x_n, x_{n-1})} - \dots = \varphi(x).$$

Примем $x = x_n = x^{(0)}$ за исходное приближение к решению вспомогательного уравнения $P(x)=0$ и построим следующие приближения по правилу $x^{(m+1)} = \varphi(x^{(m)})$ ($m=0, 1, \dots$).

Предположим, что решение с нужной точностью найдено. Его примем за следующее приближение x_{n+1} к x^* . После этого вычисляем $y_{n+1} = f(x_{n+1})$, добавляем столбец в таблицу (4.6.1) и переходим к нахождению x_{n+2} . Для этого интерполируем f по значениям в узлах $x_{n+1}, x_n, x_{n-1}, \dots, x_{n-k+1}$ и т. д.

С вычислительной точки зрения изложенный метод сложнее правила (4.6.3), так как требует решения уравнения $P(x)=0$ на каждом шаге, тогда как (4.6.3) дает явное выражение x_{n+1} через $x_n, x_{n-1}, \dots, x_{n-k}$ и требует выполнения простых арифметических операций.

Обратимся к выяснению закона изменения погрешностей $\varepsilon_n = x^* - x_n$, когда приближения x_n будут близкими к x^* . Для точного решения x^* выполняется равенство $f(x^*) = 0$ и так как, ввиду $P(x_{n+1}) = 0$,

$$f(x_{n+1}) = P(x_{n+1}) + R(x_{n+1}) = R(x_{n+1}),$$

будет $f(x^*) - f(x_{n+1}) = -R(x_{n+1})$. Если f вблизи x^* имеет непрерывную производную порядка $k+1$, при этом первая производная отлична от нуля, будут верны приводимые ниже вычисления:

$$f(x^*) - f(x_{n+1}) = f(x^*) - f(x^* - \varepsilon_{n+1}) = \varepsilon_{n+1} f'(x^* - \theta \varepsilon_{n+1}) \quad (0 < \theta < 1),$$

$$R(x_{n+1}) = \frac{(x_{n+1} - x_n)(x_{n+1} - x_{n-1}) \dots (x_{n+1} - x_{n-k})}{(k+1)!} f^{(k+1)}(\xi),$$

где ξ лежит на отрезке, содержащем x_{n-k}, \dots, x_{n+1} ,

$$\varepsilon_{n+1} = - \frac{f^{(k+1)}(\xi)}{(k+1)! f'(x^* - \theta \varepsilon_{n+1})} (\varepsilon_n - \varepsilon_{n+1})(\varepsilon_{n-1} - \varepsilon_{n+1}) \dots (\varepsilon_{n-k} - \varepsilon_{n+1}).$$

Отсюда видно, что когда приближения x_{n-k}, \dots, x_{n+1} будут близки к x^* , то ε_{n+1} будет малой величиной более высокого порядка малости, чем каждая из погрешностей $\varepsilon_n, \varepsilon_{n-1}, \dots, \varepsilon_{n-k}$, и каждую из скобок $(\varepsilon_n - \varepsilon_{n+1}), (\varepsilon_{n-1} - \varepsilon_{n+1}), \dots$ можно приближенно заменить соответственно на $\varepsilon_n, \varepsilon_{n-1}, \dots$.

Кроме того, отношение $\frac{f^{(k+1)}(\xi)}{f'(x^* - \theta \varepsilon_{n+1})}$ можно также заменить, очевидно, на $\frac{f^{(k+1)}(x^*)}{f'(x^*)}$. Приближенное равенство

$$\varepsilon_{n+1} \approx - \frac{1}{(k+1)!} \frac{f^{(k+1)}(x^*)}{f'(x^*)} \varepsilon_n \varepsilon_{n-1} \dots \varepsilon_{n-k}$$

дает, по-видимому, достаточно верное описание закона изменения погрешности ε_n при увеличении n .

§ 4.7. ИНТЕРПОЛИРОВАНИЕ С КРАТНЫМИ УЗЛАМИ

4.7.1. Существование и единственность интерполирующего многочлена. Остаток

Предположим, что на отрезке $\langle a, b \rangle$ даны m различных узлов интерполирования. Рассмотрим функцию f и будем считать, что в точке x_1 известны значения, как самой функции $f(x_1)$, так и ее производных $f'(x_1), f''(x_1), \dots, f^{(\alpha_1-1)}(x_1)$, в точке x_2 даны значения $f(x_2), f'(x_2), \dots, f^{(\alpha_2-1)}(x_2)$ и т. д. Числа $\alpha_1, \alpha_2, \dots, \alpha_m$ называются кратностями соответствующих узлов. Общее число всех исходных данных о функции f обозначим $n+1$: $\alpha_1 + \alpha_2 + \dots + \alpha_m = n+1$.

Поставим задачей найти многочлен

$$P(x) = a_0 x^n + a_1 x^{n-1} + \dots + a_n, \quad (4.7.1)$$

степени не больше n , удовлетворяющий условиям

$$P^{(i)}(x_k) = f^{(i)}(x_k) \quad (i=0, 1, \dots, \alpha_k-1, k=1, 2, \dots, m). \quad (4.7.2)$$

Эти условия дадут для определения коэффициентов a_k ($k=0, 1, \dots, n$) многочлена систему линейных уравнений. Чтобы убедиться в существовании и единственности решения системы, достаточно показать, что однородная система

$$P^{(i)}(x_k) = 0 \quad (i=0, 1, \dots, \alpha_k-1, k=1, 2, \dots, m)$$

имеет только нулевое решение. Но такая система для многочлена $P(x)$ говорит о том, что узлы x_1, x_2, \dots, x_m должны быть корнями $P(x)$ кратностей не меньше соответственно $\alpha_1, \alpha_2, \dots, \alpha_m$. Сумма кратностей корней $P(x)$ должна быть больше или равна $\alpha_1 + \alpha_2 + \dots + \alpha_m = n+1$. Но степень $P(x)$ не выше n , и иметь сумму кратностей, большую n , многочлен P может только в том случае, когда он тождественно равен нулю. Тогда все его коэффициенты a_k равны нулю и однородная система имеет, следовательно, только нулевое решение.

Таким образом, рассматриваемая интерполяционная задача с кратными узлами разрешима и имеет только одно решение, каковы бы ни были значения $f^{(i)}(x_k)$ в условиях (4.7.2).

Для многочлена $P(x)$ можно легко выписать явное выражение через узлы x_k и значения $f^{(i)}(x_k)$ при помощи определителей. Но такое представление $P(x)$ имеет сложное строение, и мы оставим его в стороне. В следующем пункте будет построено более простое представление $P(x)$. Его мы получим, воспользовавшись тем, что представление $P(x)$ зависит только от чисел $f^{(i)}(x_k)$ и не зависит от того, будет ли f аналитической функцией или не будет. Для аналитической же функции можно указать компактное выражение для P , из которого нужно нам представление P следует сравнительно просто.

Рассмотрим остаток интерполирования $R(x) = f(x) - P(x)$. Он является $n+1$ раз непрерывно дифференцируемой функцией на $\langle a, b \rangle$, удовлетворяющей условиям

$$R^{(i)}(x_k) = 0 \quad (i=0, 1, \dots, \alpha_k-1, k=1, \dots, m),$$

которые говорят о том, что узлы x_1, \dots, x_m для остатка R будут корнями, кратности которых будут не меньше соответственно $\alpha_1, \dots, \alpha_m$.

Построим для $R(x)$ одно из простейших известных представлений.

Введем многочлен степени $n+1$, связанный с узлами x_1, \dots, x_m и их кратностями $\alpha_1, \dots, \alpha_m$,

$$A(x) = (x-x_1)^{\alpha_1} \dots (x-x_m)^{\alpha_m}. \quad (4.7.3)$$

Теорема 1. Если узлы x_k ($k=1, \dots, m$) и точка интерполирования x принадлежат отрезку $[\alpha, \beta]$ и функция f имеет на этом отрезке непрерывную производную порядка $n+1$, то на $[\alpha, \beta]$ существует такая точка ξ , что для остатка интерполирования $R(x)$ верно равенство

$$R(x) = \frac{A(x)}{(n+1)!} f^{(n+1)}(\xi). \quad (4.7.4)$$

Доказательство. Чтобы отличить в обозначении точку интерполирования x от аргумента, назовем последний z и рассмотрим вспомогательную функцию

$$F(z) = f(z) - P(z) - \frac{A(z)}{A(x)} [f(x) - P(x)].$$

Она имеет на $[\alpha, \beta]$ непрерывную производную порядка $n+1$ и точки $z=x_1, \dots, z=x_m, z=x$ для нее будут нулями кратностей не ниже $\alpha_1, \dots, \alpha_m, 1$. Как обычно, мы считаем, что x не совпадает ни с одним из узлов x_k . Первая производная $F'(z)$ будет иметь, ввиду теоремы Ролля, внутри каждого отрезка между смежными точками x_1, \dots, x_m, x по меньшей мере один корень. Число таких отрезков равно m . Кроме того, узлы x_1, \dots, x_m будут корнями $F'(z)$ кратностей не меньше, чем $\alpha_1-1, \dots, \alpha_m-1$. Поэтому $F'(z)$ на $[\alpha, \beta]$ имеет нулей не меньше $(\alpha_1-1) + \dots + (\alpha_m-1) + m = \alpha_1 + \dots + \alpha_m = n+1$.

Повторив такие же рассуждения для $F'(z)$, можно прийти к заключению, что $F''(z)$ имеет по меньшей мере n нулей и т. д. и производная $F^{(n+1)}(z)$ порядка $n+1$ будет иметь на $[\alpha, \beta]$ по меньшей мере один нуль.

Назовем этот нуль ξ , и так как

$$F^{(n+1)}(z) = f^{(n+1)}(z) - \frac{(n+1)!}{A(x)} [f(x) - P(x)] = f^{(n+1)}(z) - \frac{(n+1)!}{A(x)} R(x),$$

то должно быть $f^{(n+1)}(\xi) - \frac{(n+1)!}{A(x)} R(x) = 0$, откуда следует (4.7.4).

**4.7.2. Представление $R(x)$ в случае аналитической функции f .
Формула Эрмита для многочлена $P(x)$**

Для дальнейшего изложения будет удобно изменить обозначения, указав в них явно функцию f , и вместо $P(x)$ и $R(x)$ употреблять знаки $P(f; x)$ и $R(f; x)$.

Многочлен $R(f; x)$ может быть, очевидно, записан в форме

$$P(f; x) = \sum_{k=1}^m \sum_{i=0}^{\alpha_k-1} L_{k,i}(x) f^{(i)}(x_k), \quad (4.7.5)$$

где $L_{ki}(x)$ есть многочлены степени n , зависящие от узлов x_k и их кратностей α_k . Явное выражение их дано ниже в (4.7.8).

Будем считать $f(z)$ аналитической функцией комплексной переменной z , регулярной в ограниченной замкнутой области B , содержащей внутри себя x_1, \dots, x_m и x . Контур l области будем считать спрямляемой линией.

Формула Коши $f(x) = \frac{1}{2\pi i} \int_l \frac{f(z)}{z-x} dz$, дающая представление f всюду внутри l , позволяет привести интерполяционный многочлен $P(f; x)$ для f к многочлену $P\left(\frac{1}{z-x}; x\right)$ для элементарной функции $\frac{1}{z-x}$.

Сейчас нам удобнее рассматривать не P , а остаток R :

$$\begin{aligned} R\left(\frac{1}{z-x}; x\right) &= \frac{1}{z-x} - P\left(\frac{1}{z-x}; x\right) = \\ &= \frac{1}{z-x} - \sum_{k=1}^m \sum_{i=0}^{\alpha_k-1} L_{ki}(x) \frac{i!}{(z-x_k)^{i+1}}. \end{aligned} \quad (4.7.6)$$

Рассмотрим зависимость $R\left(\frac{1}{z-x}; x\right)$ от z . Это есть правильная рациональная дробь. Для последующего полезно заметить, что $z=x$ есть простой полюс остатка с вычетом, равным 1. Общий знаменатель правой части (4.7.6) есть $(z-x)A(z)$, и после приведения к нему получим $R\left(\frac{1}{z-x}; x\right)$ в виде правильной дроби

$$R\left(\frac{1}{z-x}; x\right) = \frac{B(z, x)}{(z-x)A(z)}.$$

Здесь $B(z, x)$ есть многочлен от z , степень которого не выше $n+1$. Мы покажем сейчас, что $B(z, x)$ имеет нулевую степень относительно z и равняется $A(x)$.

Когда $|z|$ имеет большое значение, то $\frac{1}{z-x}$ разлагается в степенной ряд

$$\frac{1}{z-x} = \sum_{v=0}^{\infty} \frac{x^v}{z^{v+1}}.$$

Ввиду линейности оператора R

$$R\left(\frac{1}{z-x}; x\right) = \sum_{v=0}^{\infty} z^{-v-1} R(x^v; x).$$

Остаток интерполирования целых степеней x от нулевой до n равен нулю: $R(x^v; x) = 0$ ($v=0, 1, \dots, n$). Поэтому в предыдущем разложении все члены с $v \leq n$ равны нулю и разложение начинается с $v=n+1$

$$R\left(\frac{1}{z-x}; x\right) = \sum_{v=n+1}^{\infty} z^{-v-1} R(x^v; x).$$

Отсюда видно, что при $z \rightarrow \infty$ остаток $R\left(\frac{1}{z-x}; x\right)$ должен убывать не медленнее, чем z^{-n-2} , и в $\frac{B(z, x)}{(z-x)A(z)}$ степень числителя $B(z, x)$ относительно z должна быть на $n+2$ единицы меньше степени знаменателя. Но степень $(z-x)A(z)$ равна $n+2$, и степень $B(z, x)$ должна равняться нулю: $B(z, x) = B(x)$. Наконец, в полюсе $z=x$ вычет $R\left(\frac{1}{z-x}; x\right)$ равен единице и, стало быть, $B(x) = A(x)$. Таким образом,

$$R\left(\frac{1}{z-x}; x\right) = \frac{A(x)}{(z-x)A(z)}.$$

Отсюда и из формулы Коши для $R(f; x)$ вытекает равенство

$$R(f; x) = \frac{1}{2\pi i} \int_i R\left(\frac{1}{z-x}; x\right) f(z) dz = \frac{A(x)}{2\pi i} \int_i \frac{f(z)}{(z-x)A(z)} dz. \quad (4.7.7)$$

Вычисляя последний интеграл при помощи вычетов, можно просто найти представление $P(f; x) = f(x) - R(f; x)$.

Вычет $\frac{A(x)f(z)}{(z-x)A(z)}$ в простом полюсе $z=x$ равен $f(x)$. Найдем вычет этой функции в полюсе $z=x_k$. Когда z близко к x_k , верны приводимые ниже разложения по степеням $z-x_k$

$$f(z) = \sum_{s=0}^{\infty} \frac{1}{s!} f^{(s)}(x_k) (z-x_k)^s,$$

$$\frac{1}{z-x} = \frac{1}{(z-x_k) - (x-x_k)} = - \sum_{s=0}^{\infty} \frac{(z-x_k)^s}{(x-x_k)^{s+1}},$$

$$\frac{(z-x_k)^{\alpha_k}}{A(z)} = \sum_{s=0}^{\infty} c_s^{(k)} (z-x_k)^s.$$

Вычет функции

$$\frac{f(z)}{(z-x)A(z)} = \frac{1}{(z-x_k)^{\alpha_k}} \frac{(z-x_k)^{\alpha_k}}{A(z)} \frac{f(z)}{z-x}$$

в полюсе $z=x_k$ получится, если перемножить три приведенных выше ряда и подсчитать коэффициент при $(z-x_k)^{\alpha_k-1}$. Он равен

$$- \sum_{i=0}^{\alpha_k-1} f^{(i)}(x_k) \frac{1}{i!} \sum_{s=0}^{\alpha_k-1-i} c_s^{(k)} (x-x_k)^{-\alpha_k+s+i}.$$

Все эти вычисления приводят к следующему эрмитову представлению интерполяционного многочлена:

$$P(f; x) = \sum_{k=1}^m \sum_{i=0}^{\alpha_k-1} f^{(i)}(x_k) \frac{1}{i!} \frac{A(x)}{(x-x_k)^{\alpha_k}} \sum_{s=0}^{\alpha_k-1-i} c_s^{(k)} (x-x_k)^{i+s}. \quad (4.7.8)$$

Мы получили его, предполагая функцию f аналитической. Но в него входят только величины

$$f^{(i)}(x_k) \quad (i=0, 1, \dots, \alpha_k-1, k=1, 2, \dots, m)$$

и оно остается верным для всякой функции f , имеющей конечные значения этих величин.

Рассмотрим пример интерполирования, когда все кратности равны двум: $\alpha_k = 2$ ($k = 1, 2, \dots, m$). Степень многочлена $P(f; x)$ равна $2m - 1$, и условия интерполирования следующие:

$$P(f; x_k) = f(x_k), \quad P'(f; x_k) = f'(x_k) \quad (k = 1, 2, \dots, m). \quad (4.7.9)$$

В этом случае говорят об интерполировании соприкосновения первого порядка.

Равенство (4.7.8) в рассматриваемой задаче имеет вид

$$P(f; x) = \sum_{k=1}^m \frac{\omega^2(x)}{(x-x_k)^2 [\omega'(x_k)]^2} \times \\ \times \left\{ \left[1 - \frac{\omega''(x_k)}{\omega'(x_k)} (x-x_k) \right] f(x_k) + (x-x_k) f'(x_k) \right\}. \quad (4.7.10)$$

В справедливости его можно убедиться при помощи проверки выполнения условий (4.7.9) или при помощи составления (4.7.8) указанным выше путем.

§ 4.8. СХОДИМОСТЬ ИНТЕРПОЛЯЦИОННЫХ ПРОЦЕССОВ

При изучении сходимости интерполяционных процессов мы ограничимся наиболее простой задачей, когда интерполирование выполняется по значениям функции.

Рассмотрим интерполяционный процесс, определяемый треугольной таблицей узлов

$$X = \left\{ \begin{array}{cccc} x_1^{(1)} & & & \\ x_1^{(2)} & x_2^{(2)} & & \\ \cdot & \cdot & \cdot & \cdot \\ x_1^{(n)} & x_2^{(n)} & \dots & x_n^{(n)} \\ \cdot & \cdot & \cdot & \cdot \end{array} \right\}. \quad (4.8.1)$$

Отрезок интерполирования $[a, b]$ предполагается конечным и узлы $x_k^{(n)}$ ($k = 1, 2, \dots, n$) — лежащими на этом отрезке. Будем считать, что узлы перенумерованы в порядке роста их координат: $x_k^{(n)} < x_{k+1}^{(n)}$ ($k = 1, 2, \dots, n-1$). Допустим, что на $[a, b]$ рассматривается некоторая функция $f(x)$. Возьмем узлы $x_k^{(n)}$ ($k = 1, 2, \dots, n$), лежащие в строке

номера n , и построим по ним алгебраический многочлен $P_n(x)$ степени $n-1$, удовлетворяющий условиям

$$P_n(x_k^{(n)}) = f(x_k^{(n)}) \quad (k=1, 2, \dots, n). \quad (4.8.2)$$

Он может быть записан, например, в форме Лагранжа:

$$P_n(x) = \sum_{k=1}^n \frac{\omega_n(x)}{(x-x_k^{(n)})\omega_n'(x_k^{(n)})} f(x_k^{(n)}) = \sum_{k=1}^n l_{nk}(x) f(x_k^{(n)}).$$

Нас будет интересовать преимущественно равномерная и, в отдельных случаях, поточечная сходимость интерполирования на всем отрезке $[a, b]$:

$$\lim_{n \rightarrow \infty} P_n(x) = f(x) \quad (a \leq x \leq b). \quad (4.8.3)$$

Требуется выяснить, как между собой должны быть связаны свойства функции f и таблицы X , чтобы имело место соотношение (4.8.3) во всех точках $[a, b]$ или равномерно на $[a, b]$.

Начнем с наиболее интересного для приложений случая, когда $f(x)$ является аналитической функцией. Как мы увидим ниже, если ограничиться рассмотрением только тех случаев, которые в этой задаче можно признать «правильными», вопрос о сходимости решается при помощи сравнительно грубых признаков, содержание которых зависит от двух факторов: от области регулярности функции f и от предельной функции распределения узлов интерполирования $x_k^{(n)}$. Изложение начнем с функций распределения.

4.8.1. О предельной функции распределения узлов

Узлы интерполирования $x_k^{(n)}$ предполагаются лежащими на отрезке $[a, b]$ и, в соответствии с этим, функции распределения будут рассматриваться только на этом отрезке. Их наглядный смысл весьма прост. Вообразим, что взята единичная масса и мы ее произвольным образом распределили на $[a, b]$. Пусть x есть любая точка $[a, b]$, отличная от b . Под значением функции $\mu(x)$ в точке x будем понимать сумму масс, лежащих строго левее точки x . В частности, $\mu(a)=0$, так как левее a нет масс. В точке же $x=b$ мы положим $\mu(b)=1$.

Функция $\mu(x)$ будет обладать свойствами:

- 1) $\mu(a)=0$ при $x=a$;
- 2) $\mu(x)$ есть монотонная неубывающая функция x при $a \leq x < b$, непрерывная слева в каждой точке внутри $[a, b]$;
- 3) $\mu(b)=1$ при $x=b$.

Эти свойства мы примем за определение функции распределения и всякую функцию $\mu(x)$, обладающую ими, будем называть функцией распределения на $[a, b]$.

Допустим, что нам дана последовательность функций распределения $\mu_n(x)$ ($n=1, 2, \dots$). Говорят, что последовательность $\mu_n(x)$ сходится в основном к функции $\mu(x)$, когда сходимость имеет место во всякой точке непрерывности $\mu(x)$ внутри $[a, b]$. Если $\mu(x)$ обладает указанными выше тремя свойствами, то ее называют предельной функцией распределения для заданной последовательности $\mu_n(x)$.

Возьмем строку номера n таблицы X . Припишем каждому узлу $x_k^{(n)}$ массу $\frac{1}{n}$ и соответствующую функцию распределения обозначим $\mu_n(x)$. Это есть кусочно постоянная функция, для которой узлы $x_k^{(n)}$ будут точками разрыва со скачками $+\frac{1}{n}$.

Если существует функция распределения $\mu(x)$, к которой в основном сходится последовательность $\mu_n(x)$, то говорят, что $\mu(x)$ есть предельная функция распределения таблицы X .

В дальнейшем мы будем рассматривать только тот случай, когда таблица узлов X имеет такую функцию $\mu(x)$. Попутно отметим лишь, что если таблица X не имеет единственной предельной функции распределения, то принципиальную возможность получить картину сходимости интерполяционного процесса в этом «особом случае» дает известная теорема о возможности выбора из любой последовательности функций распределения частичной подпоследовательности μ_{n_k} ($k=1, 2, \dots$), кото-

рая сходилась бы в основном. Для нашей задачи это будет означать, что из полной последовательности интерполяционных многочленов нужно выбирать некоторые частичные последовательности и определять для каждой из них свои условия сходимости.

4.8.2. Сходимость интерполирования аналитических функций

Будем считать $f(z)$ аналитической функцией, регулярной в некоторой конечной замкнутой области B , содержащей $[a, b]$ внутри себя. Контур ее можно считать спрямляемой линией. Обозначим его l .

Как было показано в § 4.3, остаток интерполирования $r_n(z) = f(z) - P_n(z)$ может быть записан в форме контурного интеграла

$$r_n(x) = \frac{1}{2\pi i} \int_l \frac{\omega_n(x)}{\omega_n(z)} \cdot \frac{f(z)}{z-x} dz. \quad (4.8.4)$$

При определении условий стремления $r_n(x)$ к нулю большое значение имеет логарифмический потенциал

$$U(z) = \int_a^b \ln \frac{1}{|z-t|} d\mu(t). \quad (4.8.5)$$

Интеграл здесь понимается в смысле Стильтьеса.

$U(z)$ есть гармоническая функция всюду в комплексной плоскости z вне отрезка $[a, b]$. Когда z удаляется на бесконечность, $U(z)$ будет стремиться к $-\infty$. Линия уровня $U(z)=C$ при отрицательном C , большом по абсолютной величине, будет содержать внутри себя отрезок $[a, b]$ и будет сходной с окружностью большого радиуса. Обозначим ее l_C и область, лежащую внутри ее, назовем B_C .

Когда C возрастает, область B_C будет уменьшаться. Точную верхнюю границу значений C , при которых отрезок $[a, b]$ лежит внутри B_C , назовем λ и обозначим κ открытую область плоскости, где $U(z) < \lambda$. Дополнение к κ до комплексной плоскости обозначим β .

Теперь мы в состоянии сформулировать теорему о сходимости интерполирования.

Теорема 1. Если аналитическая функция $f(z)$ регулярна в некоторой области D , содержащей внутри себя β , то $r_n(x) \rightarrow 0$ при $n \rightarrow \infty$ равномерно относительно x из β .

Доказательство. Область β по условию теоремы лежит внутри D , поэтому существует такое значение $C' < \lambda$, что соответствующая ему область $B_{C'} + l_{C'}$ лежит внутри D .

Между C' и λ возьмем число C'' : $C' < C'' < \lambda$. Линия уровня $l_{C''}$ будет лежать внутри $l_{C'}$ и содержать β , а следовательно, и отрезок $[a, b]$ внутри себя.

В интегральном представлении остатка (4.8.4) за линию интегрирования примем $l_{C'}$ и будем считать, что точка x лежит на $l_{C''}$. Для остатка $r_n(x)$ верна оценка

$$|r_n(x)| \leq \frac{M}{2\pi\delta} \int_{l_{C'}} \frac{|\omega_n(x)|}{|\omega_n(z)|} ds,$$

где $M = \max_{l_{C'}} |f|$ и δ есть расстояние между линиями $l_{C'}$ и $l_{C''}$. Имеем

$$|\omega_n(z)|^{-1} = \exp \sum_{k=1}^n \ln \frac{1}{|z - x_k^{(n)}|}.$$

Каждому узлу интерполирования $x_k^{(n)}$ ($k=1, 2, \dots, n$) припишем массу $\frac{1}{n}$ и соответствующую функцию распределения обозначим $\mu_n(t)$.

$$\int_a^b \ln \frac{1}{|t-z|} d\mu_n(t) = \frac{1}{n} \sum_{k=1}^n \ln \frac{1}{|z-x_k^{(n)}|},$$

и поэтому

$$|\omega_n(x)|^{-1} = \exp n \int_a^b \ln \frac{1}{|z-t|} d\mu_n(t).$$

Точка z лежит вне отрезка $[a, b]$ и $\ln \frac{1}{|t-z|}$ является непрерывной функцией t . При неограниченном росте n функция $\mu_n(t)$ будет сходиться в основном к $\mu(t)$. По теореме Хелли о предельном переходе для интеграла Стильеса *) можно утверждать, что верно следующее соотношение:

$$\int_a^b \ln \frac{1}{|t-z|} d\mu_n(t) \rightarrow \int_a^b \ln \frac{1}{|t-z|} d\mu(t) = C' \quad (n \rightarrow \infty).$$

В интеграле z является параметром и, если проследить ход обычного доказательства теоремы Хелли, будет видно, что сходимость является равномерной относительно $z \in l_{C'}$. Поэтому существует такое n' , что при $n > n'$ и всяких $z \in l_{C'}$ будут выполняться неравенства

$$C' - \frac{1}{3} (C'' - C') < \int_a^b \ln \frac{1}{|t-z|} d\mu_n(t) < C' + \frac{1}{3} (C'' - C').$$

По аналогичным соображениям можно утверждать, что существует такое n'' , что при $n > n''$ и любых $x \in l_{C''}$ будут верны неравенства

$$C'' - \frac{1}{3} (C'' - C') < \int_a^b \ln \frac{1}{|t-x|} d\mu_n(t) < C'' + \frac{1}{3} (C'' - C').$$

Из этих двух результатов следует, что при $n > \max(n', n'')$, и всяких $z \in l_{C'}, x \in l_{C''}$ имеет место оценка

$$\begin{aligned} \int_a^b \ln \frac{1}{|z-t|} d\mu_n(t) - \int_a^b \ln \frac{1}{|x-t|} d\mu_n(t) &< \left[C' + \frac{1}{3} (C'' - C') \right] - \\ &- \left[C'' - \frac{1}{3} (C'' - C') \right] = -\frac{1}{3} (C'' - C'), \end{aligned}$$

*) В. И. Гливенко. Интеграл Стильеса, n° 14. М., 1936. И. П. Натансон, Теория функций вещественной переменной, гл. 8, § 7. М.—Л., 1950.

откуда получаются оценки для отношения $\frac{\omega_n(x)}{\omega_n(z)}$:

$$\left| \frac{\omega_n(x)}{\omega_n(z)} \right| < e^{-\frac{n}{3}(C''-C')},$$

и для остатка интерполирования $r_n(x)$:

$$|r_n(x)| < \frac{Ms}{2\pi\delta} e^{-\frac{n}{3}(C''-C')} \quad (n > \max(n', n''), x \in l_{C''}),$$

где s — длина линии $l_{C'}$.

Отсюда сразу же следует, что $r_n(x) \rightarrow 0$ ($n \rightarrow \infty$) равномерно относительно $x \in l_{C''}$, и так как $r_n(x)$ есть регулярная аналитическая функция в $l_{C''} + B_{C''}$, то $r_n(x) \rightarrow 0$ равномерно относительно x в $l_{C''} + B_{C''}$. Так как множество β лежит внутри $l_{C''} + B_{C''}$, это будет верно и для $x \in \beta$, что доказывает теорему. Из последней оценки видно также, что $r_n(x)$ стремится к нулю не медленнее, чем по закону геометрической прогрессии со знаменателем $\exp \left[-\frac{1}{3}(C''-C') \right]$.

Рассмотрим два частных случая доказанной теоремы.

1. Пусть предельная функция $\mu(x)$ соответствует равномерному распределению единичной массы на отрезке $[a, b]$:

$$\mu(x) = \frac{x-a}{b-a} \quad (a \leq x \leq b).$$

Так будет, например, в том случае, когда интерполирование выполняется по равноотстоящим узлам на $[a, b]$:

$$x_k^{(n)} = a + \frac{k-1}{n-1}(b-a) \quad (k=1, 2, \dots, n).$$

Не нарушая общности, можно считать отрезок $[a, b]$ приведенным к $[0, 1]$ и $\mu(x) = x$. Логарифмический потенциал (4.8.5) здесь есть

$$U(z) = \int_0^1 \ln \frac{1}{|t-z|} dt.$$

Ввиду

$$\int_0^1 \ln|t-z| dt = \operatorname{Re} \int_0^1 \ln(t-z) dt = \operatorname{Re} \{(1-z) \ln(1-z) + z \ln z - 1\}$$

будет

$$\begin{aligned} U(z) &= \operatorname{Re} \{1 - z \ln z - (1-z) \ln(1-z)\} = \\ &= 1 - x \ln \sqrt{x^2 + y^2} - (1-x) \ln \sqrt{(1-x)^2 + y^2} + y \operatorname{arctg} \frac{y}{x - x^2 - y^2}. \end{aligned}$$

Линии уровня потенциала $U(z)$ изображены на рис. 4.8.1. Область β ограничена линией уровня, проходящей через точки $(x=0, y=0)$ и $(x=1, y=0)$ и имеющей уравнение $U(z)=0$ или

$$x \ln \sqrt{x^2 + y^2} + (1-x) \ln \sqrt{(1-x)^2 + y^2} - y \operatorname{arctg} \frac{y}{x - x^2 - y^2} = 1.$$

Все изложенное позволяет считать доказанной приводимую ниже теорему.

Теорема 2. Пусть узлы интерполирования $x_k^{(n)}$ ($k=1, \dots, n$; $n=1, 2, \dots$) лежат на отрезке $[0, 1]$ и предельное распределение их на этом отрезке является равномерным, так что предельная функция распределения для таблицы узлов X есть $\mu(x)=x$. Если аналитическая функция $f(z)$ регулярна в замкнутой области β , изображенной на рис. 4.8.1, то интерполяционный процесс, определяемый таблицей X , сходится всюду на $[0, 1]$ равномерно относительно x .

2. Отрезок интерполирования $[a, b]$ будем считать приведенным к $[-1, 1]$. Пусть таблица X узлов $x_k^{(n)}$ имеет предельную функцию распределения

$$\mu(x) = \frac{1}{\pi} \int_{-1}^x \frac{dt}{\sqrt{1-t^2}}. \quad (4.8.6)$$

Ее часто называют функцией Чебышева.

Такой предельной функцией будет обладать, например, таблица нулей всякой системы многочленов, ортогональных на $[-1, 1]$ по весу $p(x)$, почти везде положительному на $[-1, 1]$. Это могут быть многочлены Чебышева, Лежандра, Якоби и т. д.

Рассмотрим соответствующий (4.8.6) логарифмический потенциал

$$U(z) = \frac{1}{\pi} \int_{-1}^1 \ln \frac{1}{|t-z|} \cdot \frac{dt}{\sqrt{1-t^2}} = \operatorname{Re} \left\{ \frac{1}{\pi} \int_{-1}^1 \ln \frac{1}{z-t} \cdot \frac{dt}{\sqrt{1-t^2}} \right\} = \operatorname{Re} F(z).$$

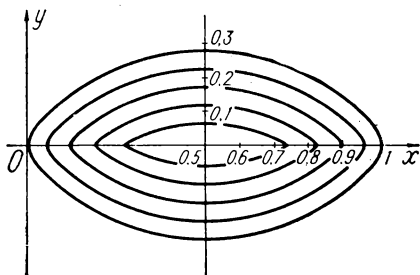


Рис. 4.8.1

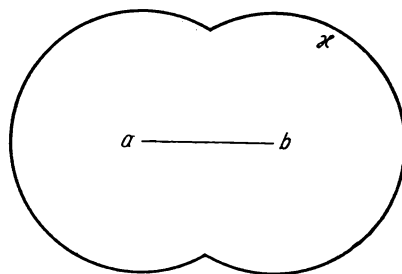


Рис. 4.8.2

В комплексной плоскости z проведем разрез вдоль действительной оси от точки $z=1$ к $-\infty$ и возьмем ту ветвь логарифма, для которой $\arg(z-t)=0$ при действительных z , больших t .

$$F'(z) = -\frac{1}{\pi} \int_{-1}^1 \frac{dt}{(z-t)\sqrt{1-t^2}}.$$

Интеграл может быть вычислен путем нахождения первообразной функции для

$$F'(z) = -\frac{1}{\sqrt{z^2-1}},$$

при этом выбирается та ветвь корня, которая имеет положительные значения при $z > 1$.

$$F(z) = \ln \frac{C}{z + \sqrt{z^2-1}}.$$

Для определения C можно воспользоваться условием, что $F(z)$ при больших положительных z имеет разложение вида

$$F(z) = \ln \frac{1}{z} + \frac{\alpha_1}{z} + \frac{\alpha_2}{z^2} + \dots$$

Получим $C=2$.

$$F(z) = \ln \frac{2}{z + \sqrt{z^2 - 1}}, \quad U(z) = \ln \frac{2}{|z + \sqrt{z^2 - 1}|}. \quad (4.8.7)$$

Линиями уровня потенциала $U(z) = \text{const} = C$ при $C < \ln 2$ будут эллипсы, имеющие общие фокусы в точках -1 и $+1$ на оси x . Для $C = \ln 2$ эллипс вырождается в отрезок $[-1, 1]$. Здесь множество β будет состоять только из прямолинейного отрезка $[-1, 1]$. Отсюда следует

Теорема 3. Пусть узлы интерполирования $x_k^{(n)}$ расположены на отрезке $[-1, 1]$ и таблица X узлов имеет предельной функцией распределения $\mu(x)$ функцию Чебышева (4.8.6). Интерполяционный процесс, определяемый такой таблицей X , сходится на отрезке $[-1, 1]$ равномерно относительно x для всякой аналитической функции $f(z)$, регулярной на отрезке $[-1, 1]$, включая его концы.

Интересно отметить, что может быть доказана теорема, обратная теореме 3.

Теорема 4. Если для таблицы X узлов интерполирования $x_k^{(n)}$, расположенных на $[-1, 1]$, интерполяционный процесс сходится во всех точках отрезка $[-1, 1]$ для всякой функции $f(x)$, аналитической на $[-1, 1]$, то для таблицы X существует предельная функция распределения узлов $x_k^{(n)}$ и это есть функция Чебышева (4.8.6).

Доказательство теоремы, к сожалению, является сложным и не может быть здесь приведено.*)

Остановимся еще на одном вопросе сходимости. В теории интерполирования известна теорема: если функция $f(z)$ — целая, то интерполяционный процесс для нее сходится равномерно на $[a, b]$, как бы ни были расположены узлы $x_k^{(n)}$ на $[a, b]$. Но легко видеть, что предположение о регулярности f всюду на комплексной плоскости, кроме бесконечности, является чрезмерно ограничительным. Утверждение о сходимости должно остаться верным, если $f(z)$ будет регулярна в некоторой конечной, но достаточно широкой около $[a, b]$, области.

Целью приводимой ниже теоремы будет указать точную наименьшую область, регулярность в которой обеспечивает сходимость интерполяционного процесса на $[a, b]$ при всякой таблице узлов $x_k^{(n)}$, взятых на отрезке $[a, b]$.

Построим два круга радиуса $b-a$ с центрами в точках a и b и обозначим κ замкнутую область, являющуюся суммой этих двух кругов (рис. 4.8.2).

*) См., например, В. И. К р ы л о в. Приближенное вычисление интегралов, гл. 12, § 2. М., 1967.

Теорема 5. Если аналитическая функция $f(z)$ регулярна в области κ , то, какова бы ни была таблица X узлов $x_k^{(n)}$, лежащих на $[a, b]$, интерполяционный процесс будет сходиться на $[a, b]$ равномерно относительно x .

Область κ является наименьшей, обеспечивающей сходимость интерполяционного процесса при любой таблице узлов $x_k^{(n)} \in [a, b]$.

Доказательство. При помощи несложных геометрических соображений можно убедиться в следующем: предположим, что x и t есть две любые точки $[a, b]$ и z — произвольная точка комплексной плоскости. Тогда, если z лежит вне κ , будет $|x-t| < |z-t|$ и $\left| \frac{x-t}{z-t} \right| < 1$, если же z принадлежит κ , то найдутся такие x и t на $[a, b]$, что будет $|x-t| \geq |z-t|$ и $\left| \frac{x-t}{z-t} \right| \geq 1$.

Так как $f(z)$ предполагается регулярной в области κ , включая границу, существует замкнутая линия l , содержащая κ внутри себя, такая, что $f(z)$ будет регулярной внутри l и на ней самой. Остаток интерполирования $r_n(x)$ представим контурным интегралом

$$r_n(x) = \frac{1}{2\pi i} \int_l \frac{\omega_n(x)}{\omega_n(z)} \frac{f(z)}{z-x} dz = \frac{1}{2\pi i} \int_l \frac{f(z)}{z-x} \prod_{k=1}^n \left(\frac{x-x_k^{(n)}}{z-x_k^{(n)}} \right) dz.$$

По сделанному выше замечанию будет $\left| \frac{x-x_k^{(n)}}{z-x_k^{(n)}} \right| < 1$, ввиду того, что x и $x_k^{(n)}$ лежат на $[a, b]$, а z лежит на l и, стало быть, вне κ . Кроме того, очевидно, существует такое число $q < 1$, что при всяких x , $x_k^{(n)} \in [a, b]$ и $z \in l$ выполняются неравенства

$$\left| \frac{x-x_k^{(n)}}{z-x_k^{(n)}} \right| \leq q < 1 \quad \text{и} \quad \left| \frac{\omega_n(x)}{\omega_n(z)} \right| \leq q^n.$$

Поэтому

$$|r_n(x)| \leq \frac{q^n}{2\pi} \max_x \int_l \frac{|f(z)|}{|z-x|} dz = Aq^n.$$

Отсюда вытекает, что $r_n(x) \rightarrow 0$ ($n \rightarrow \infty$) равномерно по $x \in [a, b]$. Проверим теперь, что область κ не может быть уменьшена. Для этого достаточно показать, что если взять любую точку $\alpha \in \kappa$, то существует такая

функция $f(z)$, регулярная в κ везде, кроме точки α , и такая система узлов $x_k^{(n)}$, что интерполяционный процесс для f будет расходиться в некоторой точке $x \in [a, b]$.

Пусть $\alpha \in \kappa$. Можно считать, что α лежит вне $[a, b]$. Возьмем функцию $f(z) = \frac{1}{z-\alpha}$. За контур l в интегральном представлении остатка $r_n(x)$ может быть принята линия, состоящая из окружности большого радиуса с центром в середине отрезка $[a, b]$, из малой окружности γ , окружающей точку α , и из двух сторон разреза, соединяющего эти окружности. Интегралы по сторонам разреза взаимно уничтожаются и для остатка получится представление

$$r_n(x) = r_n\left(\frac{1}{z-\alpha}, x\right) = \frac{\omega_n(x)}{2\pi i} \int_{\Gamma+\gamma} \frac{dz}{\omega_n(z)(z-\alpha)(z-x)}.$$

Интеграл по большой окружности Γ равен нулю, ввиду того что функция $\frac{1}{\omega_n(z)(z-\alpha)(z-x)}$ имеет в бесконечно далекой точке плоскости z нуль выше первой кратности. Интеграл по малой окружности γ берется в направлении движения часовой стрелки и равен вычету интегрируемой функции в точке $z=\alpha$, умноженному на $-2\pi i$:

$$r_n\left(\frac{1}{z-\alpha}, x\right) = \frac{\omega_n(x)}{\omega_n(\alpha)(x-\alpha)}.$$

Точка α принадлежит κ , и поэтому существуют такие точки x и t на $[a, b]$, что будет $\left|\frac{x-t}{\alpha-t}\right| \geq 1$. Точку интерполирования x закрепим, а узлы интерполирования возьмем совпадающими с точкой t : $x_k^{(n)} = t$ ($k=1, 2, \dots, n$). Это будет интерполирование с одним n -кратным узлом t . Выполняться оно будет по значению в точке t функции и производных от нее до порядка $n-1$. Интерполирующий многочлен будет отрезком разложения Тейлора для f около точки t . Многочлен $\omega_n(x)$ здесь есть $\omega_n(x) = (x-t)^n$. Для остатка интерполирования получим

$$r_n\left(\frac{1}{z-\alpha}, x\right) = \left(\frac{x-t}{\alpha-t}\right)^n \frac{1}{x-\alpha}.$$

Так как $\left|\frac{x-t}{\alpha-t}\right| \geq 1$, остаток не будет стремиться к нулю при неограниченном росте n .

Нами был построен пример расходящегося интерполяционного процесса с одним n -кратным узлом. Но, очевидно, если взять узлы $x_k^{(n)}$ различными и достаточно близкими к t и заставить их достаточно быстро приближаться к t при возрастании n , можно построить пример несходящегося интерполяционного процесса с различными узлами.

4.8.3. Некоторые вспомогательные теоремы

Мы должны перейти к изучению условий сходимости интерполяционных процессов для множеств непрерывных и непрерывно дифференцируемых на $[a, b]$ функций. Они много шире множеств аналитических функций, и условия сходимости на них должны быть более ограничительными. Такая характеристика интерполяционного процесса, как предельное распределение узлов, не может здесь быть достаточной для суждения о сходимости.*) Признаки сходимости потребуют более глубокого проникновения в существо задачи и более детального изучения всей картины вопроса. Сейчас мы рассмотрим некоторые вспомогательные факты, на которые будем опираться в исследованиях.

Лемма 1. При $x \neq 2k\pi$ ($k=0, \pm 1, \dots$) и всяких l и n ($l, n=1, 2, \dots$) верны оценки:

$$\left. \begin{aligned} |C_l^n| &= \left| \sum_{k=l}^n \cos kx \right| \leq \frac{1}{\left| \sin \frac{x}{2} \right|}, \\ |S_l^n| &= \left| \sum_{k=l}^n \sin kx \right| \leq \frac{1}{\left| \sin \frac{x}{2} \right|}. \end{aligned} \right\} \quad (4.8.8)$$

Доказательство. Рассматриваемые суммы C_l^n и S_l^n являются действительной и мнимой частью комплексной суммы

$$E_l^n = \sum_{k=l}^n e^{ikhx} = \frac{e^{i(n+1)x} - e^{ilx}}{e^{ix} - 1}.$$

Но так как

$$|E_l^n| < \frac{2}{|1 - e^{ix}|} = \frac{1}{\left| \sin \frac{x}{2} \right|},$$

должны выполняться и неравенства (4.8.8).

*) Пусть рассматривается задача об условиях сходимости интерполяционного процесса для множества всех функций, имеющих заданный порядок дифференцируемости на $[a, b]$. Такой процесс должен быть сходящимся для всякой функции, аналитической на $[a, b]$, и тогда по теореме 4 таблица X узлов интерполирования должна иметь предельную функцию распределения узлов, которая должна быть функцией Чебышева (4.8.6).

Лемма 2. При всяких x и n верна оценка

$$|\sigma_n(x)| = \left| \sum_{k=1}^n \frac{\sin kx}{k} \right| \leq 2\sqrt{\pi}. \quad (4.8.9)$$

Доказательство. σ_n является 2π -периодической и нечетной функцией. Ее достаточно рассмотреть на отрезке $[0, \pi]$ и так как при $x=0$ и $x=\pi$ $\sigma_n(x)$ обращается в нуль, то можно считать $0 < x < \pi$.

Пусть m есть целое число, для которого выполнено неравенство

$$m \leq \frac{1}{x} \sqrt{\pi} < m+1. \quad (4.8.10)$$

Тогда

$$\left| \sum_{k=1}^n \frac{\sin kx}{k} \right| \leq \left| \sum_{k=1}^m \frac{\sin kx}{k} \right| + \left| \sum_{k=m+1}^n \frac{\sin kx}{k} \right|.$$

При этом считается, что при $m=0$ справа отсутствует первая сумма, а для $m \geq n$ отсутствует вторая сумма. Ввиду $|\sin \varphi| \leq |\varphi|$, для первой суммы получим

$$\left| \sum_{k=1}^m \frac{\sin kx}{k} \right| \leq \sum_{k=1}^m \frac{kx}{k} = mx \leq \sqrt{\pi}. \quad (4.8.11)$$

Перейдем ко второй сумме. Если воспользоваться введенным в лемме 1 обозначением

$$S_l^n = \sum_{k=l}^n \sin kx, \text{ ее можно записать в виде}$$

$$\begin{aligned} \sum_{k=m+1}^n \frac{\sin kx}{k} &= \frac{S_{m+1}^{m+1}}{m+1} + \frac{S_{m+1}^{m+2} - S_{m+1}^{m+1}}{m+2} + \dots + \frac{S_{m+1}^n - S_{m+1}^{n-1}}{n} = \\ &= S_{m+1}^{m+1} \left(\frac{1}{m+1} - \frac{1}{m+2} \right) + S_{m+1}^{m+2} \left(\frac{1}{m+2} - \frac{1}{m+3} \right) + \dots + \\ &\quad + S_{m+1}^{n-1} \left(\frac{1}{n-1} - \frac{1}{n} \right) + S_{m+1}^n \frac{1}{n}. \end{aligned}$$

Отсюда, если воспользоваться второй оценкой (4.8.8),

$$\begin{aligned} \left| \sum_{k=m+1}^n \frac{\sin kx}{k} \right| &\leq \frac{1}{\sin \frac{x}{2}} \left[\left(\frac{1}{m+1} - \frac{1}{m+2} \right) + \left(\frac{1}{m+2} - \frac{1}{m+3} \right) + \dots + \right. \\ &\quad \left. + \left(\frac{1}{n-1} - \frac{1}{n} \right) + \frac{1}{n} \right] = \frac{1}{(m+1) \sin \frac{x}{2}}. \end{aligned}$$

В отношении $\frac{\sin t}{t}$ для $0 \leq t \leq \frac{\pi}{2}$ числитель возрастает медленнее знаменателя, и своего наименьшего значения отношение достигает при $t = \frac{\pi}{2}$. Стало быть, $\frac{\sin t}{t} \geq \frac{2}{\pi}$. Отсюда следует, что $\sin \frac{x}{2} \geq \frac{x}{\pi}$. Кроме того, ввиду (4.8.10), $m+1 > \frac{\sqrt{\pi}}{x}$ и

$$\left| \sum_{k=m+1}^n \frac{\sin kx}{k} \right| < \frac{1}{\frac{x}{\pi} \frac{\sqrt{\pi}}{x}} = \sqrt{\pi}.$$

Из последнего неравенства и (4.8.11) следует утверждение леммы.

Отметим полезное для дальнейшего неравенство, вытекающее из (4.8.9). Рассмотрим тригонометрические многочлены

$$\left. \begin{aligned} A(\theta) &= \frac{\cos \theta}{n-1} + \frac{\cos 2\theta}{n-2} + \dots + \frac{\cos(n-1)\theta}{1}, \\ B(\theta) &= \frac{\cos(n+1)\theta}{1} + \frac{\cos(n+2)\theta}{2} + \dots + \frac{\cos(2n-1)\theta}{n-1}. \end{aligned} \right\} \quad (4.8.12)$$

Проверим выполнение неравенства

$$|A(\theta) - B(\theta)| \leq 4\sqrt{\pi}. \quad (4.8.13)$$

В самом деле,

$$A(\theta) - B(\theta) = \sum_{k=1}^{n-1} \frac{\cos(n-k)\theta}{k} - \sum_{k=1}^{n-1} \frac{\cos(n+k)\theta}{k},$$

и так как

$$\cos(n-k)\theta - \cos(n+k)\theta = 2 \sin k\theta \sin n\theta,$$

будет

$$A(\theta) - B(\theta) = 2 \sin n\theta \sum_{k=1}^{n-1} \frac{\sin k\theta}{k}; \quad |A(\theta) - B(\theta)| \leq 2 \left| \sum_{k=1}^{n-1} \frac{\sin k\theta}{k} \right|.$$

Для проверки (4.8.13) осталось лишь воспользоваться оценкой (4.8.9).

Лемма 3. При всяких значениях $\theta_1, \theta_2, \dots, \theta_n$ ($0 \leq \theta_k \leq \pi$), различных между собой, и любых числах f_1, f_2, \dots, f_n существует четный тригонометрический многочлен $T(\theta)$, степени не больше $n-1$, выполняющий условия

$$T(\theta_i) = f_i \quad (i=1, 2, \dots, n). \quad (4.8.14)$$

Такой многочлен — единственный.

Доказательство. Всякий четный тригонометрический многочлен может быть разложен по степеням $\cos \theta$:

$$T(\theta) = a_0 + a_1 \cos \theta + a_2 \cos^2 \theta + \dots + a_{n-1} \cos^{n-1} \theta.$$

Заменим переменную θ , положив $\cos \theta = x$. При этом отрезок $0 \leq \theta \leq \pi$ изменения θ взаимно однозначно преобразуется в отрезок $[-1, 1]$ оси x . Значения θ_i перейдут в $x_i = \cos \theta_i$, также различные между собой. Тригонометрический многочлен $T(\theta)$ преобразуется в алгебраический многочлен

$$P(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1}.$$

Условия же (4.8.13) дадут для $P(x)$ систему равенств $P(x_i) = f_i$, которую можно рассматривать как условия построения многочлена степени $n-1$, интерполирующего некоторую функцию по ее значениям f_i в n различных точках x_i . Многочлен $P(x)$, как мы выяснили выше, всегда может быть построен и будет единственным. По многочлену $P(x)$, выполнив замену $x = \cos \theta$, найдем $T(\theta)$.

Лемма 4. При всяких значениях $\theta_1, \theta_2, \dots, \theta_n$ ($0 \leq \theta_k \leq \pi$; $\theta_i \neq \theta_k$, $i \neq k$) существует четный тригонометрический многочлен $T(\theta)$, степени не больше $n-1$, который выполняет неравенства

$$|T(\theta_k)| \leq 8 \sqrt[n]{\pi} \quad (k=1, 2, \dots, n) \quad (4.8.15)$$

и для которого существует на $[0, \pi]$ такая точка α , что будет

$$T(\alpha) > \ln n. \quad (4.8.16)$$

Доказательство. Рассмотрим четные тригонометрические многочлены $C_k(t)$ ($k=1, 2, \dots, n$) степени $n-1$, выполняющие условия

$$C_k(\theta_i) = \begin{cases} 0, & i \neq k; \\ 1, & i = k. \end{cases}$$

Существование таких многочленов гарантируется леммой 3. Затем введем тригонометрический многочлен

$$= A(\theta) - \sum_{k=1}^n [B(\theta_k + \theta) + B(\theta_k - \theta)] C_k(\theta),$$

где $A(\theta)$ и $B(\theta)$ указаны в (4.8.12). $U(\theta)$ есть четный тригонометрический многочлен.

Покажем, что его среднее значение равно нулю:

$$\int_0^\pi U(\theta) d\theta = \frac{1}{2} \int_{-\pi}^\pi U(\theta) d\theta = 0. \quad (4.8.17)$$

Действительно, $A(2\theta)$ есть тригонометрический многочлен без свободного члена и он, следовательно, ортогонален на $[-\pi, \pi]$ к 1. Что же касается суммы $B(\theta_k + \theta) + B(\theta_k - \theta)$, то это есть линейная комбинация, составленная из $\cos m\theta$ при $m > n$, и потому она ортогональна на $[-\pi, \pi]$ ко всякому тригонометрическому многочлену, степени не выше n . В частности, она ортогональна к многочленам $C_k(\theta)$, степень которых меньше n .

Из (4.8.17) следует, что на $[0, \pi]$ $U(\theta)$ изменяет знак и, значит, там существует такая точка α , где $U(\theta)$ обращается в нуль:

$$U(\alpha) = 0.$$

Положим

$$T(\theta) = [A(\theta + \alpha) + A(\theta - \alpha)] - \sum_{k=1}^n [B(\theta_k + \alpha) + B(\theta_k - \alpha)] C_k(\theta).$$

Это есть четный тригонометрический многочлен, степень которого не превосходит $n-1$.

Ввиду (4.8.13),

$$|T(\theta_k)| = |[A(\theta_k + \alpha) - B(\theta_k + \alpha)] + [A(\theta_k - \alpha) - B(\theta_k - \alpha)]| \leq 4\sqrt{\pi} + 4\sqrt{\pi} = 8\sqrt{\pi}$$

и для $T(\theta)$ выполняются неравенства (4.8.15).

При $\theta = \alpha$

$$T(\alpha) = A(0) + U(\alpha) = A(0) = 1 + \frac{1}{2} + \dots + \frac{1}{n-1} > \int_1^n \frac{dx}{x} = \ln n$$

и, следовательно, выполняется также (4.8.16).

Приводимая ниже лемма 5 есть простое следствие леммы 4.

Лемма 5. *Каковы бы ни были на отрезке $[a, b]$ узлы x_1, x_2, \dots, x_n ; различные между собой, существует алгебраический многочлен $P(x)$, степени не большей $n-1$, выполняющий неравенства*

$$|P(x_k)| \leq 8\sqrt{\pi} \quad (k=1, 2, \dots, n), \quad (4.8.18)$$

для которого в некоторой точке $c \in [a, b]$ будет

$$|P(c)| > \ln n. \quad (4.8.19)$$

Доказательство. Не уменьшая общности, можно считать, что отрезок $a \leq x \leq b$ линейным преобразованием приведен к $[-1, 1]$.

Выполним замену $x = \cos \theta$, переводящую $[-1, 1]$ взаимно однозначно в $[0, \pi]$. Пусть при этом точки x_k перейдут в θ_k . Последние будут различны между собой. Если

$$T(\theta) = a_0 + a_1 \cos \theta + a_2 \cos^2 \theta + \dots + a_{n-1} \cos^{n-1} \theta$$

есть многочлен, существование которого доказано в лемме 4, то

$$P(x) = T(\arccos x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{n-1} x^{n-1}$$

будет удовлетворять неравенствам (4.8.18) и (4.8.19), при этом $c = \cos \alpha$.

Перейдем теперь к доказательству важной для наших целей теоремы. Предположим, что на конечном отрезке $[a, b]$ дана бесконечная треугольная таблица X узлов $x_k^{(n)}$. Пусть на $[a, b]$ задана функция f . Возьмем в X строку номера n с узлами $x_k^{(n)}$ ($k=1, 2, \dots, n$) и построим многочлен $P_n(x)$ степени не выше $n-1$, интерполирующий f по ее значениям в узлах $x_k^{(n)}$:

$$P_n(x) = \sum_{k=1}^n \frac{\omega_n(x)}{(x-x_k^{(n)})\omega_n'(x_k^{(n)})} f(x_k^{(n)}) = \sum_{k=1}^n l_{nk}(x) f(x_k^{(n)}).$$

При исследовании условий сходимости $P_n(x)$ к $f(x)$ существенное значение имеет функция

$$\lambda_n(x) = \sum_{k=1}^n |l_{nk}(x)|.$$

Мы будем рассматривать численную величину

$$\lambda_n = \max \lambda_n(x) \quad (a \leq x \leq b). \quad (4.8.20)$$

Теорема 6. Для всякой таблицы X узлов $x_k^{(n)}$ выполняется неравенство

$$\lambda_n > \frac{\ln n}{8 \sqrt[n]{\pi}}. \quad (4.8.21)$$

Доказательство. По лемме 5 найдется многочлен $P(x)$, имеющий степень не выше $n-1$ и выполняющий неравенства (4.8.18) и (4.8.19). $P(x)$ можно представить в виде

$$P(x) = \sum_{k=1}^n l_{nk}(x) P(x_k^{(n)}).$$

Отсюда следует, ввиду $|P(x_k)| \leq 8 \sqrt[n]{\pi}$, что

$$|P(x)| \leq 8 \sqrt[n]{\pi} \sum_{k=1}^n |l_{nk}(x)|,$$

и так как $|P(c)| > \ln n$, будет

$$\lambda_n = \max_x \sum_{k=1}^n |l_{nk}(x)| > \frac{\ln n}{8 \sqrt[n]{\pi}},$$

что доказывает неравенство (4.8.21).

4.8.4. Сходимость интерполирования на множествах непрерывных и непрерывно дифференцируемых функций

Как и выше, будем считать отрезок $[a, b]$ интерполирования конечным и узлы $x_k^{(n)}$ — лежащими на этом отрезке. Пусть дана таблица узлов X (4.8.1), определяющая интерполяционный процесс.

Рассмотрим множество C_r всех функций, r раз непрерывно дифференцируемых на $[a, b]$. Нашей основной задачей будет выяснить, какими свойствами должна обладать таблица X , чтобы интерполяционный процесс сходил на $[a, b]$ для всех функций $f \in C_r$.

Главное внимание мы обратим на условия равномерной сходимости и значительно меньшее место отведем изучению поточечной сходимости.

С увеличением r множество C_r будет уменьшаться и условия сходимости будут становиться менее ограничительными. Быть может, будет интересно проследить, какое влияние на условия сходимости оказывает изменение порядка дифференцируемости r .

В некоторых случаях с целью выяснить, какое влияние оказывают на сходимость другие структурные свойства функций, мы будем отходить от множеств C_r и рассматривать, например, множества A_r функций f , имеющих на $[a, b]$ не просто непрерывную, а абсолютно непрерывную производную порядка r .

Начнем с доказательства простой теоремы, устанавливающей принципиальную возможность интерполировать равномерно и сколь угодно точно любую функцию, непрерывную на $[a, b]$.

Теорема 7. Если функция f непрерывна на $[a, b]$, то для нее существует такая таблица X , что соответствующий ей интерполяционный процесс для f будет сходиться равномерно на $[a, b]$.

Доказательство. Рассмотрим последовательность алгебраических многочленов наилучшего приближения *) к f на $[a, b]$. Она сходится к $f(x)$ равномерно на $[a, b]$. Возьмем из этой последовательности многочлен $P_{n-1}(x)$ степени $n-1$. Как известно, существует на $[a, b]$ по меньшей мере $n+1$ точек $y_1 < y_2 < \dots < y_{n+1}$, в которых разность $f - P_{n-1}$ принимает значения поочередно противоположных знаков. В каждом из промежутков (y_k, y_{k+1}) разность по меньшей мере один раз обращается в нуль и P_{n-1} принимает значение, одинаковое с f . Таких точек не меньше n . Примем n из них за узлы

интерполирования $x_k^{(n)}$ ($k=1, \dots, n$) и поместим в строку номера n таблицы X . Соответствующий интерполирующий многочлен, который выше мы обозначали $P_n(x)$, будет совпадать с $P_{n-1}(x)$. При $n \rightarrow \infty$, как отмечено несколькими строками выше, $P_n = P_{n-1}$ равномерно сходится к f .

Доказанная теорема имеет, по-видимому, лишь теоретический интерес, так как находить указанные в ней узлы $x_k^{(n)}$ крайне трудно и в нашем распоряжении нет для этого эффективных средств. Но даже если в некоторых случаях их удастся найти, то значение таблицы X с такими узлами будет, как правило, ограниченным, так как X обеспечит сходимость интерполирования для f и еще, может быть, для узкого множества функций, близких по поведению к f .

Естественно поднять вопрос, существует ли такая таблица X , которая обеспечила бы равномерную сходимость интерполирования для всякой непрерывной функции. Если бы такая таблица существовала, ее разыскание имело бы большое значение. К сожалению, как показывает теорема 8, в этой задаче необходимо дать отрицательный ответ.

Теорема 8. Какова бы ни была таблица X , существует функция f , непрерывная на $[a, b]$, для которой последовательность интерполяционных многочленов $P_n(x)$, определенных таблицей X , не будет равномерно сходиться к f при $n \rightarrow \infty$.

Доказательство. Допустим противоположное: пусть интерполяционный процесс будет равномерно сходиться для всякой функции f , непрерывной на $[a, b]$, и покажем, что это приводит к противоречию с известным фактом теории интерполирования.

Если на множестве непрерывных функций определить норму $\|f\|_C = \max_x |f|$, то оно станет полным линейным нормированным пространством, которое принято обозначать буквой C . Сходимость последовательности элементов из C означает равномерную сходимость последовательности функций.

*) Необходимые сведения о многочленах наилучшего приближения можно найти в добавлении III.

Интерполяционный многочлен

$$P_n(x) = P_n(f; x) = \sum_{k=1}^n l_{nk}(x) f(x_k^{(n)})$$

является линейным оператором, отображающим пространство C в себя. По допущению, последовательность операторов является сходящейся, так как для всякой непрерывной функции f $P_n(f; x) \rightarrow f$ равномерно.

По теореме же Банаха — Штейнгауза *) одним из условий сходимости последовательности линейных операторов, переводящих полное линейное нормированное пространство в пространство того же типа, является требование ограниченности норм операторов в совокупности:

$$\|P_n\|_C \leq M < \infty \quad (n=1, 2, \dots). \quad (4.8.22)$$

По определению нормы оператора,

$$\|P_n\|_C = \sup_{\|f\|=1} \max_x |P_n(x)| = \sup_{\|f\|=1} \max_x \left| \sum_{k=1}^n l_{nk}(x) f(x_k^{(n)}) \right| \leq \max_x \sum_{k=1}^n |l_{nk}(x)| = \lambda_n.$$

С другой стороны, найденная оценка для нормы $P_n(f; x)$ достигается для определяемой ниже непрерывной функции.

Пусть $\sum_{k=1}^n |l_{nk}(x)|$ достигает своего максимума при $x = \xi$:

$$\sum_{k=1}^n |l_{nk}(\xi)| = \lambda_n.$$

Определим функцию $f_1(x)$ в узлах $x_k^{(n)}$ равенством:

$$f_1(x_k^{(n)}) = \text{sign } l_{nk}(\xi).$$

В промежутках между соседними узлами считаем ее линейной и на отрезках $a \leq x \leq x_1^{(n)}$ и $x_k^{(n)} \leq x \leq 1$ — постоянной. Такая функция непрерывна на $[a, b]$ и имеет норму, равную единице: $\|f_1\| = 1$. Для нее многочлен $P_n(x)$ в точке $x = \xi$ принимает значение

$$P_n(f; \xi) = P_n(\xi) = \sum_{k=1}^n l_{nk}(\xi) \text{sign } l_{nk}(\xi) = \sum_{k=1}^n |l_{nk}(\xi)| = \lambda_n.$$

Поэтому верна цепочка неравенств:

$$\lambda_n = P_n(f, \xi) = \sum_{k=1}^n l_{nk}(\xi) f_1(x_k^{(n)}) \leq \sup_{\|f\|=1} \max_x \left| \sum_{k=1}^n l_{nk}(x) f(x_k^{(n)}) \right| = \|P_n\|_C.$$

*) Добавление I, § 2, теорема 2. В § 2 добавления можно найти сведения и понятия, нужные для понимания дальнейшего изложения доказательства теоремы.

Сличение двух *) полученных цепочек неравенств дает значение для нормы $P_n(f; x)$:

$$\|P_n\|_C = \lambda_n.$$

Неравенство (4.8.22) говорит о том, что должно существовать такое число M , что при всяких $n=1, 2, \dots$, должно выполняться неравенство

$$\lambda_n \leq M < \infty.$$

Но этот результат противоречит неравенству (4.8.21), доказанному в теореме 6.

Как показала теорема 8, множество всех непрерывных на $[a, b]$ функций является настолько широким, что невозможно построить таблицу X , которая обеспечила бы равномерную сходимость интерполирования для всякой функции этого множества.

Таблицы X , обладающие этим свойством, могут существовать лишь в более узких множествах функций. В этой связи, быть может, представляет интерес показать, что для множества абсолютно непрерывных на $[a, b]$ функций существует таблица X такая, что соответствующий ей интерполяционный процесс будет равномерно сходиться для любой функции множества.

Установим сначала признак такой сходимости.

Теорема 9. Для того чтобы интерполяционный процесс, определяемый таблицей X , равномерно сходилась для всякой функции f , абсолютно непрерывной на $[a, b]$, необходимо и достаточно существование такого числа $M < \infty$, чтобы частичные суммы лагранжевых коэффициентов

$$\lambda_{nj}(x) = \sum_{k=j}^n l_{nk}(x)$$

выполняли неравенство

$$|\lambda_{nj}(x)| = \left| \sum_{k=j}^n l_{nk}(x) \right| \leq M \quad (n=1, 2, \dots; j=1, 2, \dots, n). \quad (4.8.23)$$

Доказательство. Характерным представлением абсолютно непрерывных функций является следующая формула:

$$f(x) = C + \int_a^x F(t) dt = C + \int_a^b F(t) E(x-t) dt. \quad (4.8.24)$$

Здесь C есть постоянная, равная значению f в точке a : $C = f(a)$, F — любая суммируемая функция и $E(x)$ определена формулой (4.5.6).

Интерполяционный многочлен для таких функций представим в форме

$$P_n(x) = P_n(f; x) = C + \int_a^b F(t) \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) dt. \quad (4.8.25)$$

Для нас полезно найти значение ядра интеграла во всех точках отрезка $[a, b]$ и составить ясное представление о его поведении. Приведенная ниже таблица дает необходимые нам сведения:

*) Для доказательства теоремы 8 достаточным является неравенство $\lambda_n \leq \|P_n\|_C$. Неравенство же $\|P_n\|_C \leq \lambda_n$ было получено для нахождения имеющего самостоятельный интерес значения нормы оператора $P_n(f, x)$.

$$K_n(x, t) = \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) = \begin{cases} \sum_{k=1}^n l_{nk}(x) = 1, & a \leq t < x_1^{(n)}; \\ \sum_{k=2}^n l_{nk}(x) + \frac{1}{2} l_{n1}(x), & t = x_1^{(n)}; \\ \sum_{k=2}^n l_{nk}(x), & x_1^{(n)} < t < x_2^{(n)}; \\ \sum_{k=3}^n l_{nk}(x) + \frac{1}{2} l_{n2}(x), & t = x_2^{(n)}; \\ \sum_{k=3}^n l_{nk}(x), & x_2^{(n)} < t < x_3^{(n)}; \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0, & x_n^{(n)} < t \leq b. \end{cases}$$

Ядро K_n есть кусочно постоянная функция t с точками разрыва $x_k^{(n)}$. Величина скачка в $x_k^{(n)}$ равна $-l_{nk}(x)$ и значение ядра в месте разрыва есть полусумма левого и правого предельных значений:

$$K_n(x, x_k^{(n)}) = \frac{1}{2} [K_n(x, x_k^{(n)} + 0) + K_n(x, x_k^{(n)} - 0)].$$

При каждом фиксированном значении x $|K_n|$, как функция от t , достигает своего максимума на некотором открытом отрезке, ограниченном двумя какими-то смежными узлами $x_i^{(n)}$ и $x_{i+1}^{(n)}$. Что же касается зависимости K_n от x , то при каждом значении t ядро K_n есть многочлен степени $n-1$ от x .

Сходимость последовательности многочленов $P_n(x)$ равносильна сходимости последовательности линейных интегральных операторов

$$A_n F = \int_a^b F(t) \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) dt = P_n(x) - C, \quad (4.8.26)$$

преобразующих множество суммируемых на $[a, b]$ функций F в множество функций, непрерывных на $[a, b]$.

Если на множестве F ввести L -норму, положив

$$\|F\|_L = \int_a^b |F(t)| dt,$$

оно станет пространством банахова типа.

Аналогично, если на множестве непрерывных на $[a, b]$ функций f ввести C -норму

$$\|f\|_C = \max_x |f(x)|,$$

это множество станет пространством банахова типа.

Операторы $A_n F$ выполняют преобразование $L \rightarrow C$ и к сходимости последовательности $A_n F$ к $f(x) - C$ можно применить теорему Банаха — Штейнгауза.*) Необходимым и достаточным условием сходимости является выполнение двух требований.

1. Сходимость $A_n F \rightarrow f(x) - C$ должна иметь место на всюду плотном в L множестве элементов. За такое множество может быть принято множество алгебраических многочленов.**) Если F есть многочлен степени $p-1$, то соответствующая ему функция f , определенная равенством (4.8.24), есть многочлен степени p . При $n > p$ интерполирующий ее многочлен $P_n(x)$ будет совпадать с f и интерполяционный процесс для f будет, очевидно, сходящимся к f . Но в таком случае будет сходиться к $f(x) - C$ последовательность значений операторов $A_n F = P_n(x) - C$ и первое условие теоремы Банаха — Штейнгауза будет выполняться.

2. Нормы операторов A_n должны быть ограничены в совокупности.

Выясним содержание этого условия в нашей задаче. Для этого вычислим сначала C -норму $A_n F$: ***)

$$\begin{aligned} \|A_n F\|_C &= \max_x \left| \int_a^b F(t) \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) dt \right| \leq \\ &\leq \max_x \max_t \left| \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) \right| \cdot \int_a^b |F(t)| dt \end{aligned}$$

и, значит,

$$\|A_n\|_C \leq \max_{x, t} \left| \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) \right|.$$

Как мы покажем сейчас, в найденной оценке нормы оператора A_n должен иметь место знак равенства. В самом деле, пусть участвующий в оценке максимум по x и t достигается в точке $x = x_0$, $t = t_0$ и пусть t_0 принадлежит интервалу $x_i^{(n)} < t < x_{i+1}^{(n)}$. Обозначим длину этого интервала буквой l и построим функцию F_i , равную $\frac{1}{l}$ на указанном интервале и нулю вне его. Очевидно, $\|F_i\|_L = 1$. Теперь запишем для нормы оператора A_n цепочку простых неравенств. Для упрощения записи воспользуемся обозначением $K(x, t)$ для ядра, а не знаком суммы:

$$\begin{aligned} \|A_n\|_C &= \sup_{\|F\|_L \leq 1} \max_x \left| \int_a^b F(t) K(x, t) dt \right| \geq \left| \int_a^b K(x_0, t) F_i(t) dt \right| = \\ &= \left| \int_{x_i^{(n)}}^{x_{i+1}^{(n)}} \frac{1}{l} K(x_0, t_0) dt \right| = |K(x_0, t_0)| = \max_{x, t} |K(x, t)|, \end{aligned}$$

*) Добавление I, § 2, теорема 2'.

**) Плотность множества алгебраических многочленов в множестве суммируемых функций становится очевидной, если принять во внимание, что к каждой суммируемой функции F можно в метрике L приблизиться сколь угодно точно при помощи функции φ , непрерывной на $[a, b]$, к непрерывной же функции можно приблизиться сколь угодно точно и равномерно на $[a, b]$ с помощью многочлена. (См. Л. А. Люстерник, В. И. Соболев. Элементы функционального анализа, гл. I, § 8. М.—Л., 1951).

***) При вычислении C -нормы $A_n F$ мы без подробных объяснений пользовались знаками наибольших значений. Из сделанных выше пояснений зависимости ядра K_n от x и t следует, что все максимумы, участвующие в оценках нормы, существуют.

$$\|A_n\|_C \geq \max_{x, t} \left| \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) \right|.$$

Сравнение этого результата с предыдущей оценкой показывает, что для нормы A_n верно равенство

$$\|A_n\|_C = \max_{x, t} \left| \sum_{k=1}^n l_{nk}(x) E(x_k^{(n)} - t) \right|.$$

Условием сходимости операторов A_n поэтому является существование такого числа M , чтобы при всяких $n=1, 2, \dots$ выполнялось неравенство

$$\|A_n\|_C = \max_{x, t} \left| \sum_{k=1}^n l_{nk} E(x_k^{(n)} - t) \right| \leq M.$$

Но оно равносильно (4.8.23), и можно сказать, что необходимым и достаточным условием сходимости A_n , а следовательно, и равномерной сходимости интерполяционных многочленов

$$P_n(x) = P_n(f; x) = C + A_n F$$

является выполнение (4.8.23).

Теперь приведем пример такой таблицы X , что отвечающий ей интерполяционный процесс будет сходиться для всякой абсолютно непрерывной функции.

Пусть отрезок интерполирования есть $[-1, 1]$ и за узлы $x_k^{(n)}$ ($k=1, \dots, n$) приняты корни многочлена Чебышева $T_n(x) = \cos(n \arccos x)$:

$$x_k^{(n)} = \cos \frac{2(n-k)+1}{2n} \pi \quad (k=1, 2, \dots, n). \quad (4.8.27)$$

Многочлен $P_n(x)$, интерполирующий f по ее значениям в $x_k^{(n)}$, есть

$$P_n(x) = \sum_{k=1}^n l_{nk}(x) f(x_k^{(n)}), \quad l_{nk}(x) = \frac{T_n(x)}{(x - x_k^{(n)}) T_n'(x_k^{(n)})}.$$

Теорема 10. Если $f(x)$ абсолютно непрерывна на $[-1, 1]$, то при $n \rightarrow \infty$ многочлен $P_n(x)$ сходится к $f(x)$ равномерно относительно x на отрезке $[-1, 1]$.

Доказательство. По теореме 9 достаточно установить ограниченность в совокупности частичных сумм лагранжевых коэффициентов:

$$|\lambda_{nj}(x)| = \left| \sum_{k=j}^n l_{nk}(x) \right| \leq M \quad (-1 \leq x \leq 1; n=1, 2, \dots; j=1, 2, \dots, n). \quad (4.8.28)$$

В теореме 6 показано, что сумма абсолютных значений $\sum_{k=1}^n |l_{nk}(x)|$ не может быть ограниченной при $n=1, 2, \dots$, так как

$$\lambda_n = \max_x \sum_{k=1}^n |l_{nk}(x)| \geq \frac{\ln n}{8 \sqrt{\pi}}$$

для любой таблицы X узлов $x_k^{(n)}$, и λ_n^* ; следовательно, неограниченно возрастает при $n \rightarrow \infty$. Ограниченность частичных сумм $\lambda_{nj}(x)$ может выполняться только за счет того, что $l_{nk}(x)$ являются знакопеременными и при сложении их происходит уничтожение главных частей $l_{nk}(x)$. При больших n проследить за такими сокращениями в суммах

$\lambda_{nj}(x) = \sum_{k=j}^n l_{nk}(x)$ было бы сложно, но, как мы сейчас увидим, суммирование здесь может быть выполнено при помощи контурного интеграла и это очень упростит исследование сумм $\lambda_{nj}(x)$.

В комплексной плоскости z возьмем замкнутую линию Γ_j^* , содержащую внутри себя узлы $x_k^{(n)}$ ($k=j, j+1, \dots, n$) и оставляющую $x_k^{(n)}$ ($k=1, \dots, j-1$) вне себя. Тогда

$$\begin{aligned}\lambda_{nj}(x) &= \sum_{k=j}^n l_{nk}(x) = (2\pi i)^{-1} \int_{\Gamma_j} \left[1 - \frac{T_n(x)}{T_n(z)} \right] \frac{dz}{z-x} = \\ &= (2\pi i)^{-1} \int_{\Gamma_j} \frac{T_n(z) - T_n(x)}{T_n(z)} \frac{dz}{z-x}.\end{aligned}$$

В самом деле, интеграл равен сумме вычетов интегрируемой функции в особых точках, лежащих внутри Γ_j . Такими точками являются корни $x_k^{(n)}$ делителя $T_n(z)$, отвечающие $k=j, \dots, n$. По правилам, известным из теории функций комплексной переменной, вычет в точке $z = x_k^{(n)}$ легко вычисляется и равен значению в этой точке функции

$$\frac{T_n(z) - T_n(x)}{T_n'(z)} \frac{1}{z-x},$$

а так как $T_n(x_k^{(n)}) = 0$, вычет равен

$$-\frac{T_n(x)}{T_n'(x_k^{(n)}) (x_k^{(n)} - x)} = l_{nk}(x).$$

Этим доказано приведенное выше равенство для $\lambda_{nj}(x)$. Интеграл

$$(2\pi i)^{-1} \int_{\Gamma_j} \frac{dz}{z-x}$$

равен нулю или единице в зависимости от того, лежит ли x вне или внутри Γ_j . Кроме того, $|T_n(x)| \leq 1$ ($-1 \leq x \leq 1$) и для проверки (4.8.28) достаточно доказать неравенство

$$\left| \int_{\Gamma_j} \frac{dz}{T_n(z)(z-x)} \right| \leq N < \infty \quad (-1 \leq x \leq 1; n=1, 2, \dots; j=1, 2, \dots, n). \quad (4.8.29)$$

Функция $\frac{1}{T_n(z)(z-x)}$ при $z \rightarrow \infty$ убывает не медленнее, чем z^{-2} . Это дает возможность в качестве линии Γ_j взять ветвь гиперболы

$$x' = \frac{1}{2} (\rho + \rho^{-1}) \cos \theta, \quad y' = \frac{1}{2} (\rho - \rho^{-1}) \sin \theta, \quad \infty > \rho > 0, \quad \theta = \frac{j-1}{n} \pi,$$

проходимой в направлении убывания ρ . Она пересекает ось x в точке $\cos \theta$ и оставляет корни $x_k^{(n)}$ ($k = j, j+1, \dots, n$) слева от себя.

Найдем на гиперболе значения всех величин, входящих в интеграл. В теории многочленов Чебышева известно следующее представление $T_n(z)$, верное на всей комплексной плоскости:

$$T_n(z) = \frac{1}{2} [(z + \sqrt{z^2 - 1})^n + (z - \sqrt{z^2 - 1})^n].$$

$$\text{На избранной гиперболе: } z = x' + iy' = \frac{1}{2} (\rho e^{i\theta} + \rho^{-1} e^{-i\theta}),$$

$$\sqrt{z^2 - 1} = \frac{1}{2} (\rho e^{i\theta} - \rho^{-1} e^{-i\theta}), \quad z + \sqrt{z^2 - 1} = \rho e^{i\theta}, \quad z - \sqrt{z^2 - 1} = \rho^{-1} e^{-i\theta},$$

$$T_n(z) = \frac{1}{2} (\rho^n e^{in\theta} + \rho^{-n} e^{-in\theta}) = \frac{1}{2} (\rho^n + \rho^{-n}) \cos(j-1)\pi = \frac{(-1)^{j-1}}{2} (\rho^n + \rho^{-n}).$$

Дальше, ввиду того, что движение по гиперболе отвечает направлению убывания ρ , будет

$$dz = -\frac{1}{2} [\rho e^{i\theta} - \rho^{-1} e^{-i\theta}] \frac{d\rho}{\rho},$$

$$\int_{\Gamma_j} \frac{dz}{T_n(z)(z-x)} = (-1)^{j-1} \int_0^\infty \frac{\rho e^{i\theta} - \rho^{-1} e^{-i\theta}}{(\rho^n + \rho^{-n})(z-x)} \cdot \frac{d\rho}{\rho}. \quad (4.8.30)$$

Интеграл по полуоси $(0, \infty)$ разложим по схеме $\int_0^\infty = \int_0^1 + \int_1^\infty$ и в интеграле по отрезку $[0, 1]$ заменим переменную интегрирования, положив $\rho = \frac{1}{\rho'}$ ($\infty > \rho' > 1$). При замене примем во внимание, что точки гиперболы, отвечающие значениям ρ и $\frac{1}{\rho}$, располагаются симметрично относительно действительной оси x' и им отвечают сопряженные значения z . После замены, отбросив штрих у ρ' , получим

$$\begin{aligned} \int_0^\infty \frac{\rho e^{i\theta} - \rho^{-1} e^{-i\theta}}{(\rho^n + \rho^{-n})(z-x)} \cdot \frac{d\rho}{\rho} &= \int_1^\infty \frac{1}{(\rho^n + \rho^{-n})} \left[\frac{\rho e^{i\theta} - \rho^{-1} e^{-i\theta}}{z-x} + \frac{\rho^{-1} e^{i\theta} - \rho e^{-i\theta}}{\bar{z}-x} \right] \frac{d\rho}{\rho} = \\ &= 2i \int_1^\infty \frac{d\rho}{\rho(\rho^n + \rho^{-n})} \operatorname{Im} \left[\frac{\rho e^{i\theta} - \rho^{-1} e^{-i\theta}}{z-x} \right]. \end{aligned}$$

После подстановки в последний интеграл вместо z значения $z = \frac{1}{2} (\rho e^{i\theta} + \rho^{-1} e^{-i\theta})$

и несложных вычислений коэффициента при мнимой части величины, стоящей в прямоугольных скобках, найдем

$$\int_{\Gamma_j} \frac{dz}{(z-x)T_n(z)} = 4i(-1)^{j-1} \sin \theta \int_1^\infty \frac{\cos \theta - \frac{1}{2}x(\rho + \rho^{-1})}{(\rho^n + \rho^{-n})|z-x|^2} \frac{d\rho}{\rho}$$

Так как при $\rho \geq 1$ $\rho^n + \rho^{-n} \geq 2$, для наших целей достаточно доказать существование числа Q , для которого выполняется неравенство

$$\sin \theta \int_1^\infty \frac{|\cos \theta - \frac{1}{2}x(\rho + \rho^{-1})|}{|z-x|^2} \frac{d\rho}{\rho} \leq Q \quad (4.8.31)$$

$$(-1 \leq x \leq 1; n=1, 2, \dots; j=1, 2, \dots, n).$$

Положив $\frac{1}{2}(\rho + \rho^{-1}) = \xi$, придадим последнему интегралу форму

$$\sin \theta \int_1^\infty \frac{|\cos \theta - x\xi|}{\xi^2 - 2x\xi \cos \theta + x^2 - \sin^2 \theta} \frac{d\xi}{\sqrt{\xi^2 - 1}}. \quad (4.8.32)$$

Для

$$\sin \theta \frac{\cos \theta - x\xi}{\xi^2 - 2x\xi \cos \theta + x^2 - \sin^2 \theta} \frac{d\xi}{\sqrt{\xi^2 - 1}}$$

первообразной функцией является $\operatorname{arctg} \frac{\cos \theta - x\xi}{\sin \theta \sqrt{\xi^2 - 1}}$. Каждая ветвь ее ограничена на отрезке $1 \leq \xi < \infty$ некоторым числом, не зависящим от n, j, x .

Если мы хотим вычислить интеграл (4.8.32) при помощи первообразной, воспользовавшись известной связью между нею и определенным интегралом, нужно отрезок интегрирования $[1, \infty)$ разделить на участки, где $\cos \theta - x\xi$ сохраняет знак, найти приращение первообразной на каждом участке и затем сложить их, приписав приращениям знак $+$ или $-$. Так как $\cos \theta - x\xi$ есть линейная функция от ξ , таких участков будет не более двух. Ввиду же ограниченности первообразной, наверное существует число Q , при котором выполняется (4.8.31).

Обратимся теперь к задаче установления признаков сходимости интерполирования на множествах непрерывно дифференцируемых функций.

Пусть рассматривается интерполяционный процесс, определенный таблицей узлов X (4.8.1):

$$P_n(f; x) = P_n(x) = \sum_{k=1}^n \frac{\omega_n(x)}{(x - x_k^{(n)})\omega_n'(x_k^{(n)})} f(x_k^{(n)}) = \sum_{k=1}^n l_{nk}(x) f(x_k^{(n)}), \quad (4.8.33)$$

Введем функцию

$$F_{n0}(t) = F_{n0}(t; x, x_k^{(n)}) = \sum_{k=1}^n l_{nk}(x) E(t - x_k^{(n)}),$$

зависящую от узлов $x_k^{(n)}$ и положения точки интерполирования x и аналогичную ядру K_n , встречающемуся в доказательстве теоремы 9. Кроме того, нам потребуются первообразные для нее

$$F_{ns}(t; x) = F_{ns}(t) = \sum_{k=1}^n l_{nk}(x) E(t - x_k^{(n)}) \frac{1}{s!} (t - x_k^{(n)})^s;$$

определяемые начальными условиями $F_{ns}^{(j)}(a) = 0$ ($j=0, 1, \dots, s-1$).

Говорят, что функция f принадлежит классу C_r $[a, b]$, если она имеет производную $f^{(r)}$ порядка r , непрерывную на $[a, b]$.

Характерное представление таких функций дается формулой Тейлора, и мы возьмем ее в виде:

$$\begin{aligned} f(x) &= \sum_{i=0}^{r-1} c_i (x-b)^i + \int_b^x g(t) \frac{(x-t)^{r-1}}{(r-1)!} dt = \Pi_{r-1}(x) + (-1)^r \int_x^b g(t) \frac{(t-x)^{r-1}}{(r-1)!} dt = \\ &= \Pi_{r-1}(x) + (-1)^r \int_a^b g(t) E(t-x) \frac{(t-x)^{r-1}}{(r-1)!} dt, \end{aligned} \quad (4.8.34)$$

$$c_i = \frac{1}{i!} f^{(i)}(b), \quad g(t) = f^{(r)}(t).$$

Параметрами представления являются числа c_i ($i=0, 1, \dots, r-1$) и непрерывная на $[a, b]$ функция g .

Теорема 11. Для того чтобы интерполяционный процесс сходиллся при $n \rightarrow \infty$ равномерно на $[a, b]$ для всякой функции $f \in C_r$ $[a, b]$ ($r \geq 1$), необходимо и достаточно существование числа $M < \infty$, для которого выполняется неравенство

$$\int_a^b |F_{n, r-1}(t)| dt = \text{Var } F_{nr}(t) \leq M \quad (n=1, 2, \dots; a \leq x \leq b). \quad (4.8.35)$$

Доказательство. Можно считать $n \geq r$. Внесем в $P_n(x)$ (4.8.33) вместо f ее представление (4.8.34). При этом многочлен $\Pi_{r-1}(x)$ интерполируется точно:

$$\begin{aligned} P_n(x) &= \Pi_{r-1}(x) + (-1)^r \int_a^b g(t) \sum_{k=1}^n l_{nk}(x) E(t - x_k^{(n)}) \frac{1}{(r-1)!} (t - x_k^{(n)})^{r-1} dt = \\ &= \Pi_{r-1}(x) + (-1)^r \int_a^b g(t) F_{n, r-1}(t) dt. \end{aligned}$$

Сходимость многочлена $P_n(z)$ к $f(x)$ и сходимость интегрального оператора

$$A_n g = \int_a^b g(t) F_{n, r-1}(t) dt = (-1)^r [P_n(x) - \Pi_{r-1}(x)]$$

к $(-1)^r [f(x) - \Pi_{r-1}(x)]$ равносильны. A_n есть линейный оператор, определенный на множестве непрерывных на $[a, b]$ функций g . Значения оператора принадлежат тому же множеству. Введем на этом множестве норму C , полагая $\|g\| = \max_x |g(x)|$. После этого

мы можем сказать, что оператор A будет осуществлять отображение банахова пространства $C[a, b]$ в себя. Для нахождения условий сходимости мы можем, как выше, применить теорему Банаха — Штейнгауза (добавление I, § 2, теорема 2').

Рассмотрим множество алгебраических многочленов. Во-первых, оно плотно в $C[a, b]$. Во-вторых, если g есть многочлен некоторой степени m , то f есть многочлен степени $m+r-1$. При $n \geq m+r$ интерполирование f будет точным и P_n совпадает с f . Интерполяционный процесс станет стационарным и будет равномерно сходиться: $P_n \rightarrow f$. Но тогда будет сходиться равномерно последовательность

$$A_n g = P_n(x) - \Pi_{r-1}(x)$$

к

$$(-1)^r [f(x) - \Pi_{r-1}(x)] = \frac{1}{(r-1)!} \int_a^b g(t) E(t-x) (t-x)^{r-1} dt.$$

Первое условие теоремы о сходимости последовательности $A_n g$ на всюду плотном множестве элементов в $C[a, b]$ здесь, очевидно, выполняется.

Остановимся на втором условии теоремы об ограниченности в совокупности норм операторов A_n . Вычислим норму A_n :

$$\begin{aligned} \|A_n g\| &= \max_x \left| \int_a^b g(t) F_{n, r-1}(t) dt \right| \leq \max_x \left\{ \int_a^b |F_{n, r-1}(t)| dt \times \right. \\ &\quad \left. \times \max_t |g(t)| \right\} = \max_x \int_a^b |F_{n, r-1}(t)| dt \cdot \|g\|. \end{aligned}$$

Для нормы оператора A_n отсюда следует оценка сверху

$$\|A_n\| \leq \max_x \int_a^b |F_{n, r-1}(t)| dt. \quad (4.8.36)$$

Пусть максимум интеграла, стоящий в неравенстве справа, достигается при $x = x_0$:

$$\begin{aligned} \max_x \int_a^b |F_{n, r-1}(t)| dt &= \max_x \int_a^b |F_{n, r-1}(t, x)| dt = \int_a^b |F_{n, r-1}(t, x_0)| dt, \\ \|A_n\| &= \sup_{\|g\|=1} \max_x \left| \int_a^b F_{n, r-1}(t, x) g(t) dt \right| = \sup_{\|g\|=1} \left| \int_a^b F_{n, r-1}(t, x_0) g(t) dt \right|. \end{aligned}$$

Функция

$$F_{n, r-1}(t, x_0) = \sum_{k=1}^n l_{nk}(x_0) E(t-x_k^{(0)}) \frac{1}{s!} (t-x_k^{(n)})^s$$

внутри каждого из интервалов $(a, x_0^{(n)})$, $(x_1^{(n)}, x_2^{(n)})$, ..., $(x_n^{(n)}, b)$ есть некоторый многочлен от t и имеет, следовательно, либо конечное число нулей, либо является тождественно

венным нулем и функция $\text{sign } F_{n, r-1}(t, x_0)$ может на отрезке $a \leq t \leq b$ иметь только конечное число точек разрыва. Но тогда ясно, что при всяком $\varepsilon > 0$ существует такая непрерывная функция $g_\varepsilon(t)$ ($\|g_\varepsilon(t)\| \leq 1$), что будет

$$\begin{aligned} \int_a^b F_{n, r-1}(t, x_0) g_\varepsilon(t) dt &> \int_a^b F_{n, r-1}(t, x_0) \text{sign} |F_{n, r-1}(t, x_0)| dt - \varepsilon = \\ &= \int_a^b |F_{n, r-1}(t, x_0)| dt - \varepsilon = \max_x \int_a^b |F_{n, r-1}(t, x)| dt - \varepsilon. \end{aligned}$$

Поэтому

$$\begin{aligned} \|A_n\| &= \sup_{\|g\|=1} \left| \int_a^b F_{n, r-1}(t, x_0) g(t) dt \right| \geq \int_a^b F_{n, r-1}(t, x_0) g_\varepsilon(t) dt > \\ &> \max_x \int_a^b |F_{n, r-1}(t, x)| dt - \varepsilon \end{aligned}$$

и так как неравенство верно при всяком $\varepsilon > 0$, то должно быть

$$\|A_n\| \geq \max_x \int_a^b |F_{n, r-1}(t, x)| dt.$$

Сравнение с (4.8.36) позволяет сказать, что

$$\|A_n\| = \max_x \int_a^b |F_{n, r-1}(t, x)| dt.$$

Требование ограниченности норм операторов A_n в совокупности равносильно неравенству (4.8.35) и выполнение (4.8.35) является, следовательно, необходимым и достаточным условием сходимости последовательности линейных операторов A_n к предельному линейному оператору

$$A_n g \rightarrow \frac{1}{(r-1)!} \int_a^b g(t) E(t-x) (t-x)^{r-1} dt = (-1)^r [f(x) - \Pi_{r-1}(x)].$$

Отсюда следует равномерная сходимость $P_n(x)$ к $f(x)$.

Отметим частный случай доказанной теоремы, который имеет также самостоятельный интерес. Пусть $r=1$ и мы рассматриваем множество непрерывно дифференцируемых функций. Составим для этого случая условие (4.8.35):

$$\int_a^b |F_{n0}(t)| dt = \int_a^b \left| \sum_{k=1}^n l_{nk}(x) E(t-x_k^{(n)}) \right| dt$$

и, так как на отрезках между узлами $x_k^{(n)}$ интегрируемая функция имеет приводимые ниже значения

$$\sum_{k=1}^n l_{nk}(x) E(t-x_k^{(n)}) = \begin{cases} 0, & a \leq t < x_1^{(n)}; \\ l_{n1}(x), & x_1^{(n)} < t < x_2^{(n)}; \\ l_{n1}(x) + l_{n2}(x), & x_2^{(n)} < t < x_3^{(n)}; \\ \dots & \dots \\ l_{n1}(x) + l_{n2}(x) + \dots + l_{nn}(x), & x_n^{(n)} < t \leq b, \end{cases}$$

интеграл легко вычисляется и это позволяет высказать следующую теорему.

Теорема 12. Для того чтобы интерполяционный процесс, определяемый таблицей X , сходиллся равномерно на $[a, b]$ для всякой функции f , непрерывно дифференцируемой на этом отрезке, необходимо и достаточно существование числа $M < \infty$ такого, что выполняется неравенство

$$|l_{n1}(x)| (x_2^{(n)} - x_1^{(n)}) + |l_{n1}(x) + l_{n2}(x)| (x_3^{(n)} - x_2^{(n)}) + \dots + \\ + |l_{n1}(x) + \dots + l_{nn-1}(x)| (x_n^{(n)} - x_{n-1}^{(n)}) \leq M \quad (n=1, 2, \dots). \quad (4.8.37)$$

Мы приведем еще одну теорему, устанавливающую интересную связь между скоростью убывания погрешности наилучшего приближения и проблемой сходимости интерполирования.

Рассмотрим две величины, связанные с таблицей X узлов интерполирования:

$$\lambda_n(x) = \sum_{k=1}^n |l_{nk}(x)| \quad \text{и} \quad \lambda_n = \max_{x \in [a, b]} \lambda_n(x).$$

Теорема 13. Пусть f непрерывна на отрезке $[a, b]$ и E_n есть погрешность ее наилучшего приближения многочленами степени не больше n .

Если $E_{n-1} \lambda_n(x_0) \rightarrow 0$ ($n \rightarrow \infty$), то последовательность значений интерполяционных многочленов в точке x_0 стремится к значению $f(x_0)$:

$$P_n(x_0) \rightarrow f(x_0) \quad (n \rightarrow \infty).$$

Если же $E_{n-1} \lambda_n \rightarrow 0$ ($n \rightarrow \infty$), то $P_n(x)$ стремится к $f(x)$ равномерно на отрезке $[a, b]$.

Доказательство. Пусть $\Pi_{n-1}(x)$ есть многочлен степени $n-1$ наилучшего приближения к f на $[a, b]$. Интерполяционный многочлен $P_n(\Pi_{n-1}, x)$, составленный для него, будет совпадать с $\Pi_{n-1}(x)$ и поэтому

$$|P_n(f; x) - f(x)| = |[P_n(f; x) - P_n(\Pi_{n-1}, x)] + [\Pi_{n-1}(x) - f(x)]| \leq \\ \leq |P_n(f - \Pi_{n-1}, x)| + |f(x) - \Pi_{n-1}(x)|.$$

Но

$$|P_n(f - \Pi_{n-1}, x)| = \left| \sum_{k=1}^n l_{nk}(x) [f(x_k^{(n)}) - \Pi_{n-1}(x_k^{(n)})] \right| \leq E_{n-1} \lambda_n(x).$$

и

$$|f(x) - \Pi_{n-1}(x)| \leq E_{n-1},$$

значит,

$$|P_n(f; x) - f(x)| \leq [\lambda_n(x) + 1] E_{n-1}.$$

Все дальнейшее является очевидным.

Л и т е р а т у р а

1. Гончаров В. Л. Теория интерполирования и приближения функций, изд. 2-е. М., 1954.
2. Марков А. А. Исчисление конечных разностей. Одесса, 1910.
3. Натансон И. П. Конструктивная теория функций. М.—Л., 1949.
4. Стеффенсен И. Ф. Теория интерполяций. М.—Л., 1935.
5. Уиттекер Э., Робинсон Г. Математическая обработка результатов наблюдений. М.—Л., 1935.

Глава 5

ЧИСЛЕННОЕ ИНТЕГРИРОВАНИЕ

§ 5.1. КВАДРАТУРНАЯ СУММА И УСЛОВИЯ ЕЕ ПОСТРОЕНИЯ. ОСТАТОК КВАДРАТУРЫ

Здесь будет рассматриваться задача о вычислении интеграла при помощи нескольких значений интегрируемой функции. Достоинство этого метода состоит в его простоте и универсальности. Мы будем рассматривать почти исключительно задачу о вычислении простого (однократного) интеграла. Более трудную задачу о вычислении кратных интегралов мы оставим в стороне и рассмотрим только вопрос о приведении их к последовательному вычислению нескольких простых интегралов.

Начнем с проблемы вычисления определенного интеграла.

5.1.1. О квадратурной сумме

Придадим интегралу специальную форму

$$\int_a^b p(x)f(x)dx, \quad (5.1.1)$$

где $\langle a, b \rangle$ есть любой конечный или бесконечный отрезок числовой оси, $p(x)$ — некоторая фиксированная функция, которую называют весовой функцией или весом, и $f(x)$ — произвольная функция некоторого класса. Выбор такой формы для интеграла связан со следующими соображениями.

Методы вычислений, рассчитанные на очень широкие классы функций, обычно обладают невысокой точностью и, если увеличивать число значений функции, участвующих в вычислениях, показывают медленную сходимость. Поясним это простым примером. Рассмотрим интеграл в его обычной форме $\int_a^b F(x)dx$. Будем считать отрезок $[a, b]$ конечным и $F(x)$ любой интегрируемой в смысле Римана функцией. Каждый такой интеграл является пределом суммы вида $\sum_{i=1}^n F(\xi_i)\Delta x_i$, и можно, принципиально говоря, найти интеграл с любой заданной точностью, взяв достаточно

малые частичные отрезки $[x_{i-1}, x_i]$ и вычислив достаточно много значений $F(\xi_i)$.

Каждая интегральная сумма определяется способом деления $[a, b]$ на части Δx_i и выбором в каждой из них промежуточных точек ξ_i .

Когда мы ставим задачу о построении правила вычисления, одинакового для всех функций F , мы не можем отдать предпочтения одним частичным промежуткам Δx_i перед другими и вынуждены будем взять все

Δx_i одинаковыми, положив $\Delta x_i = \frac{b-a}{n} = h$.

Кроме того, на основании сходных соображений о равноправности между собой точек каждого частичного отрезка Δx_i за ξ_i мы должны будем выбрать середины частичных отрезков и принять следующее правило интегрирования:

$$\int_a^b F(x) dx \approx h \left[F\left(\frac{h}{2}\right) + F\left(\frac{3}{2}h\right) + F\left(\frac{5}{2}h\right) + \dots + F\left(b - \frac{h}{2}\right) \right].$$

Оно действительно позволяет вычислить интеграл сколь угодно точно при всякой функции F , но является весьма медленно сходящимся даже для случая аналитической функции F и требует для достижения хорошей точности нахождения интеграла весьма большого числа значений F . По этой причине в практике вычислений указанное правило применяется редко и лишь в специальных случаях. Его можно, например, применять при численном интегрировании периодических функций, где это правило, как выясняется ниже, может дать высокую точность (см. конец § 5.5).

Отметим попутно, что указанное правило, вообще говоря, становится неприменимым, если интеграл $\int_a^b F(x) dx$ является несобственным: когда $F(x)$ есть неограниченная функция или когда отрезок $[a, b]$ — бесконечный.

Лучшую точность и большее значение имеют правила численного интегрирования, рассчитанные на более узкие классы функций, которые обладают некоторыми общими свойствами. Тогда точность вычисления может быть увеличена, если заранее принять во внимание эти свойства.

Каждое правило, о котором будет говориться ниже, основано на замене интегрируемой функции на какую-либо элементарную функцию — алгебраический многочлен, рациональную функцию, тригонометрический многочлен и т. п. Чтобы такая замена имела хорошую точность, необходимо, чтобы заменяемая функция F обладала высоким порядком гладкости.

Если $F(x)$ имеет какие-нибудь особенности, мы будем заинтересованы в выделении их. Такое выделение обычно делается при помощи разложения $F(x)$ на два сомножителя $F(x) = p(x)f(x)$, где $p(x)$ имеет особенности того же типа, что и $F(x)$, а $f(x)$ есть гладкая функция. Это разложение приведет нас к интегралу вида (5.1.1). Считая вес $p(x)$ фиксированным, а $f(x)$ любой гладкой функцией, мы будем строить правило интегрирования, рассчитанное на функции, имеющие одинаковые, заранее известные особенности.

Но значение $p(x)$ не ограничивается только одним этим, а является более широким. Поясним это двумя примерами. Часто приходится вычислять несобственные интегралы вида $\int_a^\infty F(x)dx$, в которых $F(x)$ может

не иметь особенностей, быть гладкой и стремиться к нулю при $x \rightarrow \infty$. В вычислении интеграла многое зависит от того, каков закон убывания $|F(x)|$. В этом случае $F(x)$ разумно разложить на два сомножителя $F(x) = p(x)f(x)$, первый из которых $p(x)$ характеризует скорость стремления $F(x)$ к нулю, а второй $f(x)$ есть некоторая гладкая функция, допускающая хорошие приближения многочленами или рациональными функциями.

При решении граничных задач дифференциальных уравнений нередко приходится иметь дело с функциями, обращающимися в нуль на концах отрезка. Здесь естественно при интегрировании учесть это свойство, положить $F(x) = (x-a)(b-x)f(x)$ и построить правило для интегрирования с весом $p(x) = (x-a)(b-x)$.

Возвратимся к интегралу (5.1.1). Вес $p(x)$ будем считать не эквивалентным нулю и таким, что его произведение pf на любую функцию f , принадлежащую взятому множеству, абсолютно интегрируемо на $\langle a, b \rangle$.

Будем строить правила вычисления следующего вида:

$$\int_a^b p(x)f(x)dx \approx \sum_{k=1}^n A_k f(x_k). \quad (5.1.2)$$

Такое равенство часто называют формулой механических квадратур, сумму $\sum_{k=1}^n A_k f(x_k)$ — квадратурной суммой, A_k — квадратурными коэффициентами и x_k — квадратурными узлами. n , x_k и A_k ($k=1, 2, \dots, n$) являются параметрами правила (5.1.2) и их надлежит выбрать так, чтобы достигнуть «возможно лучшего» результата интегрирования для всех функций избранного класса. Заметим, что в некоторых задачах не все параметры являются произвольными. Так, например, если функция f задается таблицей значений, то мы стеснены в выборе узлов x_k : мы можем взять в качестве узлов либо все табличные значения, либо часть их.

Чтобы пояснить идею правила выбора, которой мы будем руководствоваться, достаточно ее указать для того случая, когда параметры не подчинены никаким ограничениям и выбор их свободен.

Роль числа узлов n в квадратурной сумме вполне ясна: чем больше n , тем большей точности можно достигнуть при построении правила (5.1.2). Поэтому n считают произвольным, но фиксированным числом и рассматривают задачу о выборе A_k и x_k . Правом такого выбора мы воспользуемся ниже преимущественно для одной цели — увеличения степени точности правила, понимая под этим следующее.

Рассмотрим последовательность линейно независимых функций $\omega_m(x)$ ($m=0, 1, 2, \dots$) таких, что произведения $p\omega_m$ являются абсолютно интегрируемыми. Выбор $\omega_m(x)$ подчиним условию полноты в множестве f , придав ему естественный в нашей задаче смысл. Составим линейную комбинацию

$$s_n(x) = \sum_{k=1}^n a_k \omega_k(x).$$

За «расстояние» между f и s_n примем величину

$$\rho(f, s_n) = \int_a^b |p(f-s_n)| dx.$$

Последовательность $\omega_m(x)$ условимся называть полной в множестве f , если для каждой функции f из взятого множества и всякого $\varepsilon > 0$ существует такая линейная комбинация s_n , что $\rho(f, s_n) < \varepsilon$.

Когда условие полноты выполняется, из неравенства

$$\left| \int_a^b p f dx - \int_a^b p s_n dx \right| \leq \int_a^b |p(f-s_n)| dx = \rho(f, s_n)$$

следует, что $\int_a^b p f dx$ может быть вычислен сколь угодно точно, если взять линейную комбинацию достаточно большого числа первых функций ω_m и надлежащим образом избрать численные значения a_k . Очевидно при этом, что можно достигнуть тем большей точности, чем большее значение будет иметь n .

Можно надеяться на то, что если мы при помощи выбора A_k и x_k достигнем хорошей точности в интегрировании функций ω_m , то такое правило должно дать хороший результат при интегрировании любой функции f из взятого множества. Изложенные соображения позволяют указать принцип выбора параметров A_k и x_k в правиле интегрирования.

Условимся говорить, что правило (5.1.2) имеет степень точности m , если оно дает точный результат при интегрировании $\omega_0, \omega_1, \dots, \omega_m$:

$$\int_a^b p \omega_i dx = \sum_{k=1}^n A_k \omega_i(x_k) \quad (i=0, 1, \dots, m),$$

и не точно для ω_{n+1} .

При выборе A_k и x_k ставят задачу сделать степень точности наивысшей возможной. Можно ожидать, вообще говоря, такого результата: так как число параметров A_k и x_k равно $2n$, можно надеяться правило (5.1.2) сделать точным для первых $2n$ функций ω_m и достигнуть того, чтобы степень точности правила стала равной $2n-1$. Можно предполагать также, что $2n-1$ является, вообще говоря, наивысшей возможной степенью точности. Это только гипотезы, и нужно выяснить условия, при которых они будут верными. A_k и x_k должны удовлетворять системе уравнений

$$\int_a^b p \omega_i dx = \sum_{k=1}^n A_k \omega_i(x_k) \quad (i=0, 1, \dots, 2n-1),$$

линейной относительно A_k и нелинейной относительно x_k .

За функции ω_m примем последовательность степеней x : $1, x, x^2, \dots, x^m, \dots$. Линейной комбинацией s_n здесь будет алгебраический многочлен степени n

$$s_n(x) = \sum_{k=0}^n a_k x^k = P_n(x).$$

Если отрезок $[a, b]$ конечный, то многочлены $P_n(x)$ позволяют приблизиться сколь угодно точно в равномерной метрике C к любой непрерывной на $[a, b]$ функции и будут, очевидно, обладать полнотой и в указанной выше метрике $\rho(f, s_n)$. Поэтому можно ожидать, что правила интегрирования, имеющие наивысшую алгебраическую степень точности, должны дать удовлетворительный результат для вычисления интеграла (5.1.2) при всякой непрерывной функции f .

Всюду ниже, если не будет сделано на этот счет указания, будут иметься в виду правила численного интегрирования, построенные на основе приближения f алгебраическим многочленом.

Повышение степени точности не единственная цель, которую можно поставить при построении правила (5.1.2). Можно стремиться, например, к упрощению вычислений при применении этого правила и потребовать, чтобы все коэффициенты A_k были равны между собой и правило интегрирования имело бы форму

$$\int_a^b p(x)f(x)dx \approx C [f(x_1) + f(x_2) + \dots + f(x_n)]. \quad (5.1.3)$$

В него входят $n+1$ параметров C, x_k ($k=1, 2, \dots, n$), и выбором их можно надеяться достигнуть того, чтобы правило имело степень точности не ниже n .

5.1.2. Остаток приближенной квадратуры

Остаток квадратурного правила (5.1.2)

$$R_n(f) = \int_a^b p(x)f(x)dx - \sum_{k=1}^n A_k f(x_k) = \int_a^b p f dx - Q_n(x)$$

зависит от свойств функции f и от выбора правила.

При изложении задач, которые ставятся в изучении $R(f)$, пришлось бы в значительной мере повторить то, что говорилось об остатке интерполирования. Это позволяет быть более кратким.

Две основные проблемы ставятся в исследованиях $R(f)$.

Во-первых, оценка его в наиболее часто встречающихся классах функций. Представляют интерес как грубые оценки в широких классах f , полезные в изучении вопросов сходимости, так и точные оценки в более узких классах, важные при рассмотрении зависимости $R(f)$ от свойств интегрируемой функции f . Такие оценки имеют значение также при определении числа n членов в квадратурной сумме, какое нужно взять, чтобы получить значение интеграла с требуемой точностью.

Во-вторых, исследование сходимости квадратурных процессов, т. е. выяснение условий, при которых $R_n(f) \rightarrow 0$ ($n \rightarrow \infty$). По сравнению с интерполированием это более простая задача, так как остаток квадратуры есть число и вопрос стоит о стремлении к нулю численной переменной $R_n(f)$.

Квадратурный процесс, иначе говоря, последовательность квадратурных правил, определяется двумя треугольными таблицами: таблицей узлов

$$X = \begin{pmatrix} x_1^{(1)} \\ x_1^{(2)} & x_2^{(2)} \\ x_1^{(3)} & x_2^{(3)} & x_3^{(3)} \\ \cdot & \cdot & \cdot & \cdot \end{pmatrix} \quad (5.1.4)$$

и таблицей коэффициентов

$$A = \begin{Bmatrix} A_1^{(1)} & & & & \\ A_1^{(2)} & A_2^{(2)} & & & \\ A_1^{(3)} & A_2^{(3)} & A_3^{(3)} & & \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{Bmatrix}. \quad (5.1.5)$$

В вопросах сходимости процесса приходится иметь дело с тремя факторами: с классом F функций f и таблицами X и A . Две основные задачи теории сходимости имеют следующий смысл.

1. Задан класс F функций f , и нужно определить, при каких A и X квадратурный процесс будет сходиться для всякой функции $f \in F$.

2. Задан квадратурный процесс с таблицами X и A , и нужно найти класс F функций f , для которых процесс будет сходиться.

Обе эти задачи могут быть объединены в одной более общей: нужно определить, какая должна существовать связь между классом функций F и таблицами A и X , чтобы имела место сходимость квадратурного процесса.

Некоторые из задач такого вида будут рассматриваться в параграфе о сходимости.

§ 5.2. ИНТЕРПОЛЯЦИОННЫЕ КВАДРАТУРНЫЕ ПРАВИЛА И ИХ ПОГРЕШНОСТИ

Часто возникает потребность вычислить интеграл в том случае, когда известна таблица значений функции f . Мы можем пользоваться для этого табличными значениями f_k , выбирая в качестве узлов x_k квадратурного правила лишь табличные значения аргумента, и лишены возможности произвольного выбора x_k . В распоряжении вычислителя остается право выбора коэффициентов A_k , и нашей первой задачей будет выяснить, какая степень точности может быть при этом достигнута.

Для упрощения изложения будем считать отрезок интегрирования $[a, b]$ конечным, хотя в дальнейшем для многих вопросов такое предположение не является обязательным.

Один из возможных способов построения правила интегрирования состоит в следующем. Пусть узлы x_k избраны каким-либо образом и фиксированы. Интерполируем функцию f по ее значениям $f(x_k)$ ($k=1, 2, \dots, n$) при помощи алгебраического многочлена степени $n-1$

$$f(x) = P(x) + r(x), \quad P(x) = \sum_{k=1}^n \frac{\omega(x)}{(x-x_k)\omega'(x_k)} f(x_k), \quad (5.2.1)$$

$$\omega(x) = (x-x_1) \dots (x-x_n).$$

Подстановка полученного представления f в интеграл (5.1.1) дает равенство

$$\int_a^b p(x)f(x)dx = \int_a^b p(x)P(x)dx + \int_a^b p(x)r(x)dx.$$

Если здесь отбросить интеграл с остаточным членом $r(x)$, мы получим правило приближенного вычисления интеграла, которое характеризуется определенным законом выбора коэффициентов A_k . В связи со способом получения, оно названо интерполяционным правилом:

$$\int_a^b p(x)f(x)dx \approx \sum_{k=1}^n A_k f(x_k), \quad A_k = \int_a^b p(x) \frac{\omega(x)}{(x-x_k)\omega'(x_k)} dx. \quad (5.2.2)$$

Погрешность его имеет следующее выражение через остаток интерполирования $r(x)$:

$$R_n = \int_a^b p f dx - \sum_{k=1}^n A_k f(x_k) = \int_a^b p(x)r(x)dx. \quad (5.2.3)$$

Теорема 1. Для того чтобы квадратурное правило (5.1.2) было точным для всяких алгебраических многочленов степени $n-1$, необходимо и достаточно, чтобы оно было интерполяционным.

Доказательство. Начнем с проверки необходимости. Положим

$$\tilde{f}(x) = \frac{\omega(x)}{(x-x_i)\omega'(x_i)} = \omega_i(x).$$

Это есть многочлен степени $n-1$, принимающий в узлах x_k значения $\omega_i(x_k) = 0$ ($k \neq i$) и $\omega_i(x_i) = 1$. Так как правило (5.1.2) предполагается точным для всякого многочлена степени $n-1$, оно верно и для $\omega_i(x)$ и должно быть

$$\int_a^b p(x)\omega_i(x)dx = \int_a^b p(x) \frac{\omega(x)}{(x-x_i)\omega'(x_i)} dx = \sum_{k=1}^n A_k \omega_i(x_k) = A_i$$

и правило (5.1.2) действительно является интерполяционным.

Теперь докажем достаточность. Пусть f есть произвольный многочлен, степень которого не больше $n-1$. Если f интерполировать по значениям в узлах x_k , то, в силу единственности интерполяционного многочлена, должно быть точным равенство

$$\tilde{f}(x) = \sum_{k=1}^n \frac{\omega(x)}{(x-x_k)\omega'(x_k)} f(x_k).$$

Кроме того, правило предполагается интерполяционным и, стало быть, его коэффициенты имеют значения (5.2.2). Поэтому верны равенства

$$\int_a^b p(x) f(x) dx = \sum_{k=1}^n f(x_k) \int_a^b p(x) \frac{\omega(x)}{(x-x_k)\omega'(x_k)} dx = \sum_{k=1}^n A_k f(x_k)$$

и правило (5.1.2) точно выполняется для \tilde{f} .

Доказанная теорема позволяет утверждать, что всякое квадратурное правило (5.1.2), степень точности которого не меньше $n-1$, является интерполяционным.

Рассмотрим остаток (5.2.3). Каждое известное представление погрешности интерполирования $r(x)$ порождает соответствующее ему представление остатка R_n . Например, известно, что остаток $r(x)$ для каждой функции f с конечными значениями в узлах x_k и точке x имеет форму (4.3.7), которая в нашем случае может быть записана как

$$r(x) = \omega(x) f(x_1, x_2, \dots, x_n, x).$$

Это дает возможность утверждать, что для любой f с конечными значениями на $[a, b]$ и такой, что произведение $p(x)\omega(x)f(x)$ интегрируемо на $[a, b]$, погрешность R_n интерполяционного правила (5.2.2) представима в виде

$$R_n = \int_a^b p(x) \omega(x) f(x_1, x_2, \dots, x_n, x) dx. \quad (5.2.4)$$

Если, кроме того, f имеет непрерывную производную порядка n на $[a, b]$, то

$$f(x_1, x_2, \dots, x_n, x) = \frac{1}{n!} f^{(n)}(\xi),$$

где ξ есть некоторая точка на $[a, b]$. Для остатка квадратуры R_n в этом случае верно равенство

$$R_n = \frac{1}{n!} \int_a^b p(x) \omega(x) f^{(n)}(\xi) dx. \quad (5.2.5)$$

Рассмотрим функции f , для которых производная порядка n непрерывна и ограничена по абсолютному значению числом M_n :

$$|f^{(n)}(x)| \leq M_n \quad (a \leq x \leq b). \quad (5.2.6)$$

Для них из (5.2.5) следует оценка остатка

$$|R_n| \leq \frac{M_n}{n!} \int_a^b |p(x)\omega(x)| dx. \quad (5.2.7)$$

Полученная оценка является точной, когда произведение $p(x)\omega(x)$ сохраняет знак на $[a, b]$, и достигается для многочлена

$$f(x) = \frac{x^n}{n!} M_n + a_1 x^{n-1} + \dots + a_n.$$

Оценку R_n в классе (5.2.6), точную при любых $p(x)$ и узлах x_k , можно без труда получить, если воспользоваться характерным представлением классов непрерывно дифференцируемых функций, о котором мы говорили в конце § 4.5.

Всякая функция, имеющая на $[a, b]$ непрерывную производную порядка n , представима в виде

$$f(x) = \sum_{i=0}^{n-1} \frac{(x-a)^i}{i!} f^{(i)}(a) + \int_a^b f^{(n)}(t) E(x-t) \frac{(x-t)^{n-1}}{(n-1)!} dt. \quad (5.2.8)$$

Если (5.2.8) внести в остаток общего квадратурного правила

$$R_n(f) = \int_a^b p(x)f(x) dx - \sum_{k=1}^n A_k f(x_k)$$

и изменить порядок интегрирования по переменным x и t , что при предположениях о конечности $[a, b]$ и абсолютной интегрируемости $p(x)$ является возможным, мы получим для $R_n(f)$ следующее равенство:

$$R_n(f) = \sum_{i=0}^{n-1} \frac{1}{i!} f^{(i)}(a) R_n[(x-a)^i] + \int_a^b f^{(n)}(t) K_n(t) dt, \quad (5.2.9)$$

$$K_n(t) = \int_a^b p(x) \frac{(x-t)^{n-1}}{(n-1)!} dx - \sum_{x_k > t} A_k \frac{(x_k-t)^{n-1}}{(n-1)!}$$

$$\{t \neq a, x_k \quad (k=1, 2, \dots, n)\}.$$

Если же квадратурное правило является интерполяционным и, следовательно, точным для всяких многочленов степени не выше $n-1$, то

$$R_n[(x-a)^i] = 0 \quad (i=0, 1, \dots, n-1)$$

и остаток такого правила будет иметь представление

$$R_n(f) = \int_a^b f^{(n)}(t) K_n(t) dt. \quad (5.2.10)$$

Точная его оценка в классе функций, имеющих непрерывную производную порядка n , удовлетворяющую неравенству (5.2.6), будет

$$|R_n(f)| \leq M_n \int_a^b |K_n(t)| dt. \quad (5.2.11)$$

§ 5.3. ПРАВИЛА НЬЮТОНА — КОТЕСА

Интерполяционные квадратурные правила с равноотстоящими узлами рассматривались еще Ньютоном. Котесом была составлена таблица коэффициентов A_k для них в случае постоянной весовой функции $p(x) = 1$ и $n=1(1)10$.

Отрезок интегрирования $[a, b]$ разделим на n одинаковых частей длины $h = \frac{b-a}{n}$ и точки деления $a+kh$ ($k=0, 1, \dots, n$) примем за узлы интерполяционного квадратурного правила. Само правило запишем в виде

$$\int_a^b p(x) f(x) dx \approx (b-a) \sum_{k=0}^n B_k^n f(a+kh), \quad (5.3.1)$$

$$B_k^n = (b-a)^{-1} A_k = (b-a)^{-1} \int_a^b p(x) \frac{\omega(x)}{(x-a-kh)\omega'(a+kh)} dx,$$

$$\omega(x) = (x-a)(x-a-h) \dots (x-a-nh).$$

Коэффициентам B_k^n можно придать другую форму, если ввести новую переменную t , положив $x = a+th$ ($0 \leq t \leq n$):

$$\begin{aligned} \omega(x) &= h^{n+1} t(t-1)(t-2) \dots (t-n), \quad x-a-kh = h(t-k), \\ \omega'(a+kh) &= (-1)^{n-k} h^n k! (n-k)!, \end{aligned}$$

$$B_k^n = \frac{(-1)^{n-k}}{nk!(n-k)!} \int_0^n p(a+th) \frac{t(t-1)\dots(t-n)}{t-k} dt. \quad (5.3.2)$$

Остановимся более подробно на случае постоянной весовой функции $p(x) \equiv 1$.

$$\int_a^b f(x) dx \approx (b-a) \sum_{k=0}^n B_k^n f(a+kh), \quad (5.3.3)$$

$$B_k^n = \frac{(-1)^{n-k}}{nk!(n-k)!} \int_0^n \frac{t(t-1)\dots(t-n)}{t-k} dt.$$

Как говорилось выше, Котесом были вычислены коэффициенты B_k^n для $n=1(1)10$:

$$n=1; \quad B_0^1 = B_1^1 = \frac{1}{2};$$

$$n=2; \quad B_0^2 = B_2^2 = \frac{1}{6}, \quad B_1^2 = \frac{4}{6};$$

$$n=3; \quad B_0^3 = B_3^3 = \frac{1}{8}, \quad B_1^3 = B_2^3 = \frac{3}{8};$$

$$n=4; \quad B_0^4 = B_4^4 = \frac{7}{90}, \quad B_1^4 = B_3^4 = \frac{32}{90}, \quad B_2^4 = \frac{12}{90};$$

$$n=5; \quad B_0^5 = B_5^5 = \frac{19}{288}, \quad B_1^5 = B_4^5 = \frac{75}{288}, \quad B_2^5 = B_3^5 = \frac{50}{288};$$

$$n=6; \quad B_0^6 = B_6^6 = \frac{41}{840}, \quad B_1^6 = B_5^6 = \frac{216}{840}, \quad B_2^6 = B_4^6 = \frac{27}{840},$$

$$B_3^6 = \frac{272}{840};$$

$$n=7; \quad B_0^7 = B_7^7 = \frac{751}{17280}, \quad B_1^7 = B_6^7 = \frac{3577}{17280}, \quad B_2^7 = B_5^7 = \frac{1323}{17280},$$

$$B_3^7 = B_4^7 = \frac{2989}{17280};$$

$$\begin{aligned}
n=8; \quad B_0^8=B_8^8 &= \frac{989}{28350}, \quad B_1^8=B_7^8 = \frac{5888}{28350}, \quad B_2^8=B_6^8 = \frac{-928}{28350}, \\
B_3^8=B_5^8 &= \frac{10496}{28350}, \quad B_4^8 = \frac{-4540}{28350}; \\
n=9; \quad B_0^9=B_9^9 &= \frac{2857}{89600}, \quad B_1^9=B_8^9 = \frac{15741}{89600}, \quad B_2^9=B_7^9 = \frac{1080}{89600}, \\
B_3^9=B_6^9 &= \frac{19344}{89600}, \quad B_4^9=B_5^9 = \frac{5778}{89600}; \\
n=10; \quad B_0^{10}=B_{10}^{10} &= \frac{16067}{598752}, \quad B_1^{10}=B_9^{10} = \frac{106300}{598752}, \quad B_2^{10}=B_8^{10} = \frac{-48525}{598752}, \\
B_3^{10}=B_7^{10} &= \frac{272400}{598752}, \quad B_4^{10}=B_6^{10} = \frac{-260550}{598752}, \quad B_5^{10} = \frac{427368}{598752}.
\end{aligned}$$

В настоящее время эта таблица значительно продолжена. Но так как при больших n правила Ньютона — Котеса почти не применяются в вычислениях, мы ограничимся только приведенной таблицей, достаточной для наших целей. Уже при беглом рассмотрении ее можно заметить, что изменение B_k^n при возрастании k , начиная с $n=4$ и особенно при $n \geq 6$, имеет «неправильности», которые трудно считать приемлемыми. При $n=8$ и $n=10$ эти неправильности становятся особенно очевидными, так как некоторые из коэффициентов B_k^n являются отрицательными.

Заметим, что при каждом n сумма коэффициентов B_k^n всегда равна единице. В этом легко убедиться, если в правиле (5.3.3) положить $f \equiv 1$. Поэтому появление среди B_k^n отрицательных чисел вызовет рост $\sum_{k=0}^n |B_k^n|$, что может повлечь за собой потерю точности при вычислении $\sum_{k=0}^n B_k^n f(a+kh)$, так как в сумме, как правило, будут встречаться слагаемые разных знаков. Кроме того, значения $f(a+kh)$ мы знаем, обычно, только приближенно. Если погрешности известных значений f оцениваются числом ε , то погрешность в вычислении суммы $\sum_{k=0}^n B_k^n f(a+kh)$ должна быть оценена величиной $\varepsilon \sum_k |B_k^n|$, которая может иметь большие значения, если сумма $\sum_k |B_k^n|$ будет большой. Поэтому при изучении

правила Котеса имеет значение исследовать, при каких n среди коэффициентов B_k^n встречаются отрицательные и какими будут значения B_k^n при больших n .

Теорема 1. В правиле Ньютона — Котеса (5.3.3) для всех $n \geq 10$ существуют отрицательные B_k^n .

Доказательство. Не уменьшая общности рассуждений, можно считать отрезок $[a, b]$ совпадающим с $[-1, 1]$. Рассмотрим следующую формулу с $n+1$ узлами, один из которых фиксирован в точке 1:

$$\int_{-1}^1 f(x) dx \approx \sum_{k=1}^n A_k f(y_k) + A f(1) \quad (5.3.4)$$

и одновременно с ней возьмем формулу с $m+1$ ($m < n$) узлами, из которых два фиксированы в точках -1 и 1 :

$$\int_{-1}^1 f(x) dx \approx p_{-1} f(-1) + \sum_{i=1}^{m-1} p_i f(x_i) + p_{+1} f(1). \quad (5.3.5)$$

В § 5.7 будет доказано, что узлы последнего правила и его коэффициенты могут быть выбраны так, что равенство (5.3.5) будет выполняться точно для всякого многочлена степени $2m-1$, при этом узлы x_i тогда должны быть корнями многочлена Якоби степени $m-1$ индексов 1, 1: $y(x) = P_{m-1}^{(1,1)}(x)$.

Предположим, что y_k и x_i перенумерованы в порядке возрастания.

Лемма 1. Если правило интегрирования (5.3.4) верно для многочленов степени $2m-1$ ($n > m$) и если $A_k > 0$ ($k=1, \dots, n$), то $y_n > x_{m-1}$.

Доказательство. Применим правила (5.3.4) и (5.3.5) к многочлену

$$f(x) = \frac{y^2(x)(1-x^2)}{x-x_{m-1}},$$

имеющему степень $2m-1$. Оба правила должны дать точный результат, причем второе из них дает, очевидно, нуль. После сравнения результатов получим равенство

$$\sum_{k=1}^n A_k f(y_k) = 0.$$

Ввиду $n > m$ не все слагаемые в сумме равны нулю. Среди них должны быть положительные и отрицательные. Но при $x_{m-1} < x < 1$ будет $f(x) > 0$ и при $-1 < x < x_{m-1}$ $f(x) \leq 0$. Следовательно, должно быть $y_n > x_{m-1}$.

Для оценки наибольшего корня $P_{m-1}^{(1,1)}(x)$ нам потребуется

Лемма 2. Для наибольшего корня x_{m-1} многочлена Якоби $P_{m-1}^{(1,1)}(x)$ верно неравенство

$$1 - x_{m-1} < \frac{8}{(m-1)(m+2) + 4} \quad (m \geq 3). \quad (5.3.6)$$

Чтобы не прерывать рассуждений, связанных со знаками коэффициентов формулы Котеса, отложим на несколько строк доказательство леммы и допустим, что она верна. Рассмотрим правило Ньютона — Котеса для отрезка $[-1, 1]$:

$$\int_{-1}^1 f(x) dx \approx \sum_{k=0}^n A_k f\left(-1 + \frac{2k}{n}\right) \quad (5.3.7)$$

и предположим, что его коэффициенты A_k положительны.

Пусть n есть нечетное число. Правило (5.3.7) верно для многочленов степени n , и можно считать $2m-1=n$, $m=0,5(n+1)$. Ввиду лемм 1 и 2, должно быть

$$-1 + \frac{2(n-1)}{n} > x_{m-1} > 1 - \frac{32}{(n-1)(n+5)+16},$$

откуда следует

$$\frac{1}{n} < \frac{16}{(n-1)(n+5)+16} \quad \text{и} \quad n < 11.$$

Предположим, что n есть число четное. Правило (5.3.7) верно для всех многочленов степени n , так как оно интерполяционное. Кроме того, оно верно для всякой нечетной функции, ввиду того что интеграл $\int_{-1}^1 f dx$ тогда равен нулю и сумма также равна нулю, так как узлы $y_k = -1 + \frac{2k}{n}$ расположены симметрично относительно начала координат и $A_{n-k} = A_k$. Поэтому правило (5.3.7) точно для многочленов степени $n+1$ и можно считать $2m-1=n+1$, $m=0,5(n+2)$.

Согласно с леммами 1 и 2, должны выполняться неравенства

$$-1 + \frac{2(n-1)}{n} > x_{m-1} > 1 - \frac{32}{n(n+6)+16}.$$

Отсюда

$$\frac{1}{n} < \frac{16}{n(n+6)+16} \quad \text{и} \quad n < 8.$$

Можно считать теорему 1 доказанной.

Осталось еще доказать лемму 2. Известно, что многочлены Якоби $P_n^{(\alpha, \beta)}(x)$ любых индексов α, β удовлетворяют дифференциальному уравнению

$$(1-x^2)y'' + [\beta - \alpha - (\alpha + \beta + 2)x]y' + n(n + \alpha + \beta + 1)y = 0.$$

Для $y = P_{m-1}^{(1, 1)}(x)$ это уравнение будет

$$[(1-x^2)^2 y']' + (m-1)(m+2)(1-x^2)y = 0. \quad (5.3.8)$$

Пусть в уравнение вместо y подставлен многочлен $y(x) = P_{m-1}^{(1, 1)}(x)$. Проинтегрируем обе части равенства от x_{m-1} до 1. Ниже нам придется иметь дело только с наибольшим корнем x_{m-1} , и для упрощения записи мы обозначим его x , отбросив индекс $m-1$.

$$(1-x^2)^2 y'(x) = (m-1)(m+2) \int_x^1 (1-t^2)y(t) dt =$$

$$= (m-1)(m+2) \sum_{\nu=1}^{m-1} \frac{1}{\nu!} y^{(\nu)}(x) \int_x^1 (1-t^2)(t-x)^\nu dt.$$

В наибольшем корне x все производные $y^{(\nu)}(x)$ положительны. Сохраним справа только два первых члена и отбросим остальные. Когда $m > 3$, правая часть уменьшится и получится неравенство

$$(1-x^2)^2 y'(x) \geqslant \\ \geqslant (m-1)(m+2) \left[y'(x) \int_x^1 (1-t^2)(t-x) dt + \frac{1}{2} y''(x) \int_x^1 (1-t^2)(t-x)^2 dt \right].$$

Из уравнения (5.3.8) вытекает, что в корне многочлена Якоби будет

$$y''(x) = \frac{4x}{1-x^2} y'(x).$$

Внесем это значение $y''(x)$ в последнее неравенство и сократим результат на $(1-x)^2 y'(x)$, что не равно нулю. Получится

$$(1+x)^2 \geqslant (m-1)(m+2) \left[\frac{3+x}{12} + \frac{x}{15} \frac{4+x}{1+x} \right] (1-x).$$

Дробь

$$\frac{4+x}{1+x} = 1 + \frac{3}{1+x}$$

убывает при $-1 < x \leqslant 1$, поэтому правая часть уменьшится, если ее заменить значением $\frac{5}{2}$, соответствующим $x=1$. После замены и сокращения на $1+x$, будет

$$1+x > \frac{(m-1)(m+2)}{4} (1-x).$$

Отсюда сразу следует утверждение леммы 2.

Чтобы получить представление об особенностях формулы Ньютона — Котеса при большом числе узлов, найдем асимптотическое представление B_k^n для больших n . С этой целью рассмотрим интеграл

$$I = \int_0^n \frac{x(x-1)\dots(x-n)}{x-k} dx,$$

входящий в выражение (5.3.3). Очевидно,

$$x(x-1)\dots(x-n) = \frac{\Gamma(x+1)}{\Gamma(x-n)}.$$

На основании известного равенства для функции $\Gamma(z)$:

$$\frac{1}{\Gamma(z)} = \frac{\Gamma(1-z) \sin \pi z}{\pi}$$

можно написать

$$\frac{\Gamma(x+1)}{\Gamma(x-n)} = (-1)^n \frac{\Gamma(x+1) \Gamma(n+1-x) \sin \pi x}{\pi}$$

и

$$I = (-1)^n \int_0^n \frac{\Gamma(x+1) \Gamma(n+1-x)}{\pi(x-k)} \sin \pi x \, dx.$$

Представим интеграл как сумму трех слагаемых:

$$\int_0^n = \int_0^3 + \int_3^{n-3} + \int_{n-3}^n = \alpha + \beta + \gamma.$$

Рассмотрим сначала слагаемое β . В теории функции $\Gamma(z)$ известно разложение логарифмической производной

$$\frac{\Gamma'(z)}{\Gamma(z)} = -\frac{1}{z} + C + \sum_{k=1}^{\infty} \left(\frac{1}{k} - \frac{1}{k+z} \right),$$

где C — постоянная Эйлера.* Из него следует, что при $z > 0$ отношение $\frac{\Gamma'(z)}{\Gamma(z)}$ есть монотонная возрастающая функция. Поэтому

$$\frac{\Gamma'(x+1)}{\Gamma(x+1)} - \frac{\Gamma'(n+1-x)}{\Gamma(n+1-x)}$$

будет монотонной возрастающей функцией при $-1 < x < n+1$ и, следовательно,

$$\ln \Gamma(x+1) \Gamma(n+1-x)$$

и произведение

$$\Gamma(x+1) \Gamma(n+1-x)$$

при $3 \leq x \leq n-3$ будут иметь наибольшее значение при $x = n-3$:

$$0 \leq \Gamma(x+1) \Gamma(n+1-x) \leq \Gamma(4) \Gamma(n-2) = 6\Gamma(n-2).$$

Так как при всяких x верно неравенство

$$\left| \frac{\sin \pi x}{\pi(x-k)} \right| \leq 1,$$

то

$$|\beta| \leq 6\Gamma(n-2)n = \frac{6\Gamma(n+1)}{(n-2)(n-1)} = O\left[\frac{\Gamma(n+1)}{n^2}\right].$$

* См., например, Янке и Эмде. Таблицы функций. М., 1959, стр. 108, 109.

Из двух интегралов α и γ , ввиду их равноправности, достаточно рассмотреть один, например α . Будем считать $1 \leq k \leq n-1$. В рассуждениях для нас полезными являются следующие известные факты: *) при больших значениях z верны равенства

$$\psi(z) = \frac{\Gamma'(z)}{\Gamma(z)} = \ln z + O\left(\frac{1}{z}\right); \quad \psi'(z) = O\left(\frac{1}{z}\right).$$

Воспользовавшись формулой Тейлора и предыдущими равенствами, получим:

$$\ln \Gamma(n+1-x) = \ln \Gamma(n+1) - x \frac{\Gamma'(n+1)}{\Gamma(n+1)} + O\left(\frac{1}{n}\right),$$

$$\Gamma(n+1-x) = \Gamma(n+1) e^{-x \ln n} \left[1 + O\left(\frac{1}{n}\right) \right].$$

При $0 \leq x \leq 3$ верно, очевидно, равенство

$$\Gamma(x+1) \frac{\sin \pi x}{\pi(x-k)} = -\frac{x}{k} + O\left(\frac{x^2}{k}\right).$$

Стало быть,

$$\alpha = \int_0^3 \Gamma(n+1) e^{-x \ln n} \left[1 + O\left(\frac{1}{n}\right) \right] \left[-\frac{x}{k} + O\left(\frac{x^2}{k}\right) \right] dx,$$

и так как

$$\int_0^3 e^{-x \ln n} x dx = \frac{1}{\ln^2 n} - \frac{1}{n^3} \left(\frac{3}{\ln n} + \frac{1}{\ln^2 n} \right),$$

$$\int_0^3 e^{-x \ln n} x^2 dx = \frac{2}{\ln^3 n} - \frac{1}{n^3} \left(\frac{9}{\ln n} + \frac{6}{\ln^2 n} + \frac{2}{\ln^3 n} \right),$$

то

$$\alpha = -\frac{\Gamma(n+1)}{k \ln^2 n} \left[1 + O\left(\frac{1}{\ln n}\right) \right].$$

Сходные вычисления для интеграла γ дадут

$$\gamma = (-1)^{n-1} \frac{\Gamma(n+1)}{(n-k) \ln^2 n} \left[1 + O\left(\frac{1}{\ln n}\right) \right].$$

Два последних результата и ранее полученная оценка для β позволяют построить для I асимптотическое выражение

$$I = (-1)^{n-1} \frac{\Gamma(n+1)}{\ln^2 n} \left[\frac{1}{k} + \frac{(-1)^n}{n-k} \right] \left[1 + O\left(\frac{1}{\ln n}\right) \right],$$

которое приводит к нужному нам асимптотическому представлению котесова коэффициента B_k^n :

$$B_k^n = \frac{(-1)^{k-1} n!}{k! (n-k)! n \ln^2 n} \left[\frac{1}{k} + \frac{(-1)^n}{n-k} \right] \left[1 + O\left(\frac{1}{\ln n}\right) \right] \quad (1 \leq k \leq n-1), \quad (5.3.9)$$

Аналогично, для B_0^n и B_n^n получаются равенства

$$B_0^n = B_n^n = \frac{1}{n \ln n} \left[1 + O\left(\frac{1}{\ln n}\right) \right]. \quad (5.3.10)$$

Найденные выражения для B_k^n позволяют сказать, что при больших n среди B_k^n будут как положительные, так и отрицательные, превосходящие по абсолютной величине любое наперед заданное число. Весьма часты будут случаи, когда смежные коэффициенты B_k^n и B_{k+1}^n будут иметь разные знаки. Это заставляет думать, что при больших n правила Ньютона — Котеса становятся малоприменимыми для вычислений.

§5.4. НЕКОТОРЫЕ ПРОСТЕЙШИЕ ПРАВИЛА НЬЮТОНА — КОТЕСА

В вычислениях наиболее часто употребляются правила Ньютона — Котеса с малым числом узлов. Они имеют невысокую точность, и для уменьшения погрешности отрезков интегрирования $[a, b]$ нужно предварительно разделить на достаточно большое число малых интервалов, к каждому из них применить избранное правило и затем взять сумму по всем интервалам.

5.4.1. Правило трапеций

Положим $n=1$. Интерполирование в этом случае выполняется по двум значениям $f(a)$ и $f(b)$, которые принимает f на концах отрезка $[a, b]$. Равенство (5.3.3) имеет вид

$$\int_a^b f(x) dx \approx \frac{b-a}{2} [f(a) + f(b)] \quad (5.4.1)$$

и является простейшим правилом трапеций. Ввиду $p(x) \equiv 1$ и $\omega(x) = (x-a)(x-b)$, погрешность (5.2.5) правила равна

$$R(f) = \frac{1}{2} \int_a^b (x-a)(x-b) f''(\xi) dx.$$

Так как множитель $(x-a)(x-b)$ сохраняет знак на $[a, b]$, то при

условии непрерывности второй производной от f на $[a, b]$ должна существовать такая точка η , для которой

$$R(f) = f''(\eta) \int_a^b (x-a)(x-b) dx.$$

Вычисление последнего интеграла приведет к такому выражению для погрешности (5.4.1):

$$R(f) = -\frac{(b-a)^3}{12} f''(\eta) \quad (a \leq \eta \leq b). \quad (5.4.2)$$

Разделим теперь отрезок $[a, b]$ на n одинаковых частей длины $h = \frac{1}{n}(b-a)$ и рассмотрим отрезок $[a+kh, a+(k+1)h]$. Для вычисления интеграла по нему применим равенство (5.4.1):

$$\begin{aligned} \int_{a+kh}^{a+(k+1)h} f(x) dx &= \frac{h}{2} [f_k + f_{k+1}] + R_k \quad [f_k = f(a+kh)], \\ R_k &= -\frac{h^3}{12} f''(\eta_k), \quad a+kh \leq \eta_k \leq a+(k+1)h. \end{aligned}$$

Сумма таких результатов по всем частичным отрезкам приведет к общей формуле трапеций:

$$\int_a^b f(x) dx = \frac{b-a}{n} \left[\frac{1}{2} f_0 + f_1 + \dots + f_{n-1} + \frac{1}{2} f_n \right] + R. \quad (5.4.3)$$

Здесь

$$\begin{aligned} R &= R_0 + R_1 + \dots + R_{n-1} = -\frac{h^3}{12} [f''(\eta_0) + f''(\eta_1) + \dots + f''(\eta_{n-1})] = \\ &= -\frac{(b-a)^3}{12n^2} \frac{1}{n} [f''(\eta_0) + \dots + f''(\eta_{n-1})]. \end{aligned}$$

Величина $\frac{1}{n} [f''(\eta_0) + \dots + f''(\eta_{n-1})]$ есть среднее арифметическое значение, составленное из n значений второй производной f'' в n точках отрезка $[a, b]$. Оно лежит между максимальным M и минимальным m значения-

ми f'' на $[a, b]$, а так как непрерывная функция принимает всякое значение между M и m , на $[a, b]$ существует такая точка ξ , что

$$R = -\frac{(b-a)^3}{12n^2} f''(\xi) \quad (a \leq \xi \leq b). \quad (5.4.4)$$

5.4.2. Правило парабол [формула Симпсона]

Перейдем к случаю $n=2$. Интерполирование f выполняется по трем значениям в точках a , $c = \frac{1}{2}(a+b)$, b . Квадратурное правило (5.3.3) будет

$$\int_a^b f(x) dx \approx \frac{b-a}{6} [f(a) + 4f(c) + f(b)]. \quad (5.4.5)$$

Равенство является точным для всех многочленов второй степени. Но необходимо заметить, что если f есть функция нечетная относительно точки c , являющейся серединой отрезка $[a, b]$, т. е. если $f(c-t) = -f(c+t)$ при всяких t , то левая и правая части в (5.4.5) обращаются в нуль и равенство будет выполняться точно. В частности, оно будет точным для $f = (x-c)^3$. Это позволит утверждать, что правило (5.4.5) является точным для всех многочленов третьей степени.

Чтобы найти погрешность (5.4.5), построим многочлен третьей степени $P_3(x)$, удовлетворяющий условиям

$$P_3(a) = f(a), \quad P_3(c) = f(c), \quad P_3'(c) = f'(c), \quad P_3(b) = f(b).$$

$P_3(x)$ интерполирует $f(x)$ по двум однократным узлам a и b и одному двукратному узлу c :

$$f(x) = P_3(x) + r(x).$$

Здесь $r(x)$ есть остаток интерполирования.

$$\int_a^b f(x) dx = \int_a^b P_3(x) dx + \int_a^b r(x) dx.$$

Так как (5.4.5) точно для всякого многочлена третьей степени, то

$$\int_a^b f(x) dx = \frac{b-a}{6} [P_3(a) + 4P_3(c) + P_3(b)] + \int_a^b r(x) dx =$$

$$= \frac{b-a}{6} [f(a) + 4f(c) + f(b)] + \int_a^b r(x) dx$$

и $\int_a^b r(x) dx$ есть погрешность (5.4.5):

$$R = \int_a^b r(x) dx.$$

Предположим, что f имеет на $[a, b]$ непрерывную производную четвертого порядка. Для остатка $r(x)$ интерполирования тогда верно равенство вида (4.7.4), которое в нашем случае будет таким:

$$r(x) = \frac{(x-a)(x-c)^2(x-b)}{4!} f^{IV}(\xi) \quad (a \leq \xi, x \leq b)$$

и, стало быть,

$$R = \frac{1}{24} \int_a^b (x-a)(x-c)^2(x-b) f^{IV}(\xi) dx.$$

Множитель $(x-a)(x-c)^2(x-b)$ не меняет знака на $[a, b]$. Обычное в таких случаях рассуждение показывает, что на $[a, b]$ существует точка η такая, что верно равенство

$$R = \frac{1}{24} f^{IV}(\eta) \int_a^b (x-a)(x-c)^2(x-b) dx.$$

После несложных вычислений для остатка получится

$$R = -\frac{1}{90} \left(\frac{b-a}{2} \right)^5 f^{IV}(\eta).$$

Разделим отрезок $[a, b]$ на четное число n равных частей длины $h = \frac{b-a}{n}$. Возьмем удвоенный частичный отрезок $[a + (k-1)h, a + (k+1)h]$ и применим к нему правило парабол (5.4.5) с остатком

$$\int_{a+(k-1)h}^{a+(k+1)h} f(x) dx = \frac{h}{3} [f_{k-1} + 4f_k + f_{k+1}] + R_k, \quad [f_k = f(a+kh)].$$

Применив это равенство к отрезкам $[a, a+2h]$, $[a+2h, a+4h]$, ... и сложив почленно результаты, построим общее правило парабол или правило Симпсона

$$\int_a^b f(x) dx = \frac{b-a}{3n} [f_0 + f_n + 2(f_2 + f_4 + \dots + f_{n-2}) + 4(f_1 + f_3 + \dots + f_{n-1})] + R \quad (5.4.6)$$

Остаток

$$R = -\frac{1}{90} h^5 [f^{IV}(\eta_1) + f^{IV}(\eta_3) + \dots + f^{IV}(\eta_{n-1})].$$

Так как

$$\frac{2}{n} [f^{IV}(\eta_1) + \dots + f^{IV}(\eta_{n-1})] = f^{IV}(\xi),$$

где ξ есть некоторая точка отрезка $[a, b]$, для R верно равенство

$$R = -\frac{(b-a)^5}{180n^4} f^{IV}(\xi) \quad (a \leq \xi \leq b). \quad (5.4.7)$$

5.4.3. Правило «трех восьмых»

При $n=3$ формула (5.3.3) приведет к ньютонову правилу «трех восьмых»

$$\begin{aligned} \int_a^b f(x) dx = H \left[\frac{1}{8} f(a) + \frac{3}{8} f\left(a + \frac{1}{3} H\right) + \right. \\ \left. + \frac{3}{8} f\left(a + \frac{2}{3} H\right) + \frac{1}{8} f(a+H) \right] + R, \end{aligned} \quad (5.4.8)$$

$$H = b - a.$$

Пусть n есть число, кратное трем. Вновь разделим $[a, b]$ на n равных частей $h = \frac{b-a}{n}$. Применив правило (5.4.8) к строеным отрезкам $[a, a+3h]$, $[a+3h, a+6h]$, ... и сложив результаты, получим общее правило «трех восьмых», сходное с правилом Симпсона:

$$\begin{aligned} \int_a^b f(x) dx = \frac{3}{8} \frac{b-a}{n} [(f_0 + f_n) + 2(f_3 + f_6 + \dots + f_{n-3}) + \\ + 3(f_1 + f_2 + f_4 + f_5 + \dots + f_{n-2} + f_{n-1})] + R. \end{aligned} \quad (5.4.9)$$

§ 5.5. КВАДРАТУРНЫЕ ПРАВИЛА НАИВЫСШЕЙ АЛГЕБРАИЧЕСКОЙ СТЕПЕНИ ТОЧНОСТИ

Пусть в правиле численного интегрирования

$$\int_a^b p(x)f(x)dx \approx \sum_{k=1}^n A_k f(x_k) \quad (5.5.1)$$

$\langle a, b \rangle$ есть любой конечный или бесконечный отрезок и весовая функция $p(x)$ такова, что ее произведение на любую неотрицательную степень x абсолютно интегрируемо на $\langle a, b \rangle$:

$$\int_a^b |p(x)x^i|dx < \infty \quad (i \geq 0).$$

Кроме того, как всюду выше, будем считать $p(x)$ не эквивалентной нулю:

$$\int_a^b |p(x)|dx > 0.$$

Правило при фиксированном n содержит $2n$ параметров x_k, A_k и выбрать их можно так, чтобы равенство (5.5.1) выполнялось точно для всех алгебраических многочленов степени не выше $2n-1$ или, что равносильно, чтобы выполнялись равенства

$$\int_a^b p(x)x^i dx = \sum_{k=1}^n A_k x_k^i \quad (i=0, 1, \dots, 2n-1).$$

В каком случае и каким путем это может быть достигнуто, увидим ниже.

5.5.1. Построение правила и его единственность

Выясним условия, при которых (5.5.1) точно выполняется для всех многочленов степени $2n-1$.

Нам удобнее иметь дело не с узлами x_k , а с многочленом $\omega(x) = (x-x_1)(x-x_2)\dots(x-x_n)$. Знания x_k и многочлена $\omega(x)$, очевидно, равносильны. Но если мы хотим нахождение x_k заменить нахождением $\omega(x)$, мы обязаны будем показать, что корни $\omega(x)$ действительны, различны и принадлежат отрезку $\langle a, b \rangle$.

Теорема 1. Для того чтобы правило (5.5.1) было точным для всех многочленов степени не больше $2n-1$, необходимо и достаточно выполнения условий:

1) правило (5.5.1) — интерполяционное:

$$A_k = \int_a^b p(x) \frac{\omega(x)}{(x-x_k)\omega'(x_k)} dx; \quad (5.5.2)$$

2) многочлен $\omega(x)$ ортогонален на $\langle a, b \rangle$ по весу $p(x)$ ко всякому многочлену $Q(x)$, степени меньшей n :

$$\int_a^b p(x) \omega(x) Q(x) dx = 0. \quad (5.5.3)$$

Доказательство. Необходимость первого условия очевидна: если равенство (5.5.1) верно для всякого многочлена степени меньшей $2n$, то оно верно для многочленов степени меньшей n и должно быть, по теореме 1 § 5.1, интерполяционным.

Необходимость второго условия проверяется столь же просто. Пусть $Q(x)$ — любой многочлен степени меньшей n . Положим $f(x) = \omega(x)Q(x)$. Это есть многочлен, степень которого меньше $2n$. Для него правило (5.5.1) должно быть точным. Но так как $f(x_k) = 0$ ($k=1, 2, \dots, n$), правая часть (5.5.1) есть нуль и должно выполняться равенство (5.5.3).

Докажем достаточность условий теоремы. Допустим, что f есть произвольный многочлен степени меньшей $2n$. Разделив f на ω , можно представить f в форме $f(x) = \omega(x)Q(x) + \rho(x)$, где степени $Q(x)$ и $\rho(x)$ меньше n . Кроме того, ввиду $\omega(x_k) = 0$, будет $f(x_k) = \rho(x_k)$.

$$\int_a^b p(x)f(x)dx = \int_a^b p(x)\omega(x)Q(x)dx + \int_a^b p(x)\rho(x)dx.$$

Первый интеграл правой части равен нулю по второму условию и, так как по первому условию правило (5.5.1) интерполяционное, верно равенство

$$\int_a^b p(x)\rho(x)dx = \sum_{k=1}^n A_k \rho(x_k).$$

Но $\rho(x_k) = f(x_k)$ и должно также быть верным равенство

$$\int_a^b p(x)f(x)dx = \sum_{k=1}^n A_k f(x_k),$$

что доказывает достаточность условий теоремы.

Доказанная теорема приводит вопрос о возможности построения равенства (5.5.1), точного для всяких многочленов степени меньшей $2n$, к проблеме существования многочлена $\omega(x)$, обладающего свойством ортогональности (5.5.3).

Теорема 2. Если весовая функция $p(x)$ не меняет знак на $\langle a, b \rangle$, например остается неотрицательной, то существует и при этом единственный многочлен

$$\omega(x) = x^n + a_1 x^{n-1} + \dots + a_n,$$

ортогональный на $\langle a, b \rangle$ по весу $p(x)$ ко всякому многочлену степени меньшей n .

Доказательство. Будем искать многочлен $\omega(x)$ в форме разложения по степеням x , как указано в формулировке теоремы. Для определения коэффициентов a_1, \dots, a_n условия ортогональности дадут систему n уравнений

$$\int_a^b p(x) [x^n + a_1 x^{n-1} + \dots + a_n] x^i dx = 0 \quad (i=0, 1, \dots, n-1).$$

Нам достаточно убедиться в том, что соответствующая однородная система

$$\int_a^b p(x) [a_1 x^{n-1} + \dots + a_n] x^i dx = 0 \quad (i=0, 1, \dots, n-1)$$

имеет только нулевое решение, так как отсюда следует, что определитель рассматриваемой системы отличен от нуля и она имеет единственное решение. Если выписать подробно уравнения однородной системы для $i=0, 1, \dots, n-1$, умножить их последовательно на a_n, a_{n-1}, \dots, a_1 и сложить, получится равенство

$$\int_a^b p(x) [a_1 x^{n-1} + \dots + a_n]^2 dx = 0.$$

Если бы многочлен $a_1 x^{n-1} + \dots + a_n$ не был бы тождественным нулем, он мог бы обращаться в нуль не больше чем в $n-1$ точках и равенство

не могло бы выполняться, так как $p(x) \geq 0$ и $\int_a^b p(x) dx > 0$. Значит, $a_1 x^n + \dots + a_n$ тождественно равняется нулю, все его коэффициенты a_1, \dots, a_n равны, следовательно, нулю и однородная система имеет только нулевое решение.

Теорема 3. Если $p(x)$ не меняет знак на $\langle a, b \rangle$ и многочлен $\omega(x)$ ортогонален на $\langle a, b \rangle$ по весу $p(x)$ ко всякому многочлену $Q(x)$, степени меньше n , то все корни многочлена $\omega(x)$ действительные, различные и лежат внутри $\langle a, b \rangle$.

Доказательство. Рассмотрим корни многочлена $\omega(x)$, которые лежат внутри $\langle a, b \rangle$ и имеют нечетную кратность. Пусть таких корней m и это есть $\xi_1, \xi_2, \dots, \xi_m$. Нам достаточно показать, что $m=n$, так как отсюда следует, что никаких других корней у $\omega(x)$ нет и все корни ξ_k — простые.

Допустим противоположное: $m < n$ и покажем, что это противоречит свойству ортогональности. Составим многочлен

$$\rho(x) = (x - \xi_1) \dots (x - \xi_m).$$

Его степень m меньше n , и для него должно выполняться равенство

$$\int_a^b p(x) \omega(x) \rho(x) dx = 0.$$

Но сразу же видно, что это равенство не может быть выполнено, так как $\omega(x)$ и $\rho(x)$ имеют внутри $\langle a, b \rangle$ одинаковые точки перемены знака и произведение $\omega \rho$ сохраняет знак на $\langle a, b \rangle$. Кроме того, $\omega(x) \rho(x)$ обращается в нуль только в конечном числе точек, так как ω и ρ отличны от тождественного нуля. Ввиду того что вес $p(x)$ также сохраняет знак на $\langle a, b \rangle$ и не эквивалентен нулю, интеграл $\int_a^b p \omega \rho dx$ должен быть отличен от нуля, а это противоречит предыдущему.

Во всех предшествующих рассуждениях число n могло быть любым целым и положительным.

Доказанные теоремы позволяют высказать следующее утверждение.

Если вес $p(x)$ сохраняет знак на $\langle a, b \rangle$, то квадратурное правило (5.5.1), верное для многочленов степени не выше $2n-1$, существует при всяких $n=1, 2, \dots$ и является единственным для каждого n .

Осталось выяснить, будет ли $2n-1$ наивысшей возможной степенью точности. Ответ дает

Теорема 4. Если $p(x)$ сохраняет знак на $\langle a, b \rangle$, то ни при каком выборе x_k и A_k равенство (5.5.1) не может быть верным для всех многочленов степени $2n$.

Доказательство. Для проверки правильности утверждения достаточно построить многочлен, имеющий степень $2n$, для которого (5.5.1) не может быть выполнено точно. Положим $f(x) = \omega^2(x)$. Это есть положительный многочлен степени $2n$. Для него

$$\int_a^b p f dx \neq 0,$$

а сумма

$$\sum_{k=1}^n A_k f(x_k) = \sum_{k=1}^n A_k \omega^2(x_k)$$

равна нулю, так как $\omega(x_k) = 0$ ($k = 1, 2, \dots, n$), и равенство (5.5.1) не может выполняться точно.

Отсюда следует, что при знакопостоянной весовой функции $p(x)$ степень точности $2n-1$ действительно является наивысшей возможной.

5.5.2. Два замечания о квадратурных коэффициентах

Покажем, что в правиле наивысшей алгебраической степени точности, отвечающей неотрицательной весовой функции $p(x)$, все коэффициенты A_k положительны. Это утверждение есть следствие приводимой ниже теоремы.

Теорема 5. Если $p(x) \geq 0$ и квадратурное правило (5.5.1) верно для всех многочленов степени $2n-2$, то все коэффициенты A_k в нем положительны.

Доказательство. Положим

$$f(x) = \left[\frac{\omega(x)}{x-x_i} \right]^2.$$

Это есть многочлен степени $2n-2$ и для него равенство (5.5.1) должно быть верным. Но

$$f(x_k) = \begin{cases} 0, & k \neq i, \\ \omega'^2(x_i), & k = i \end{cases}$$

и, следовательно,

$$\int_a^b p(x) \left[\frac{\omega(x)}{x-x_i} \right]^2 dx = A_i \omega'^2(x_i),$$

$$A_i = \int_a^b p(x) \left[\frac{\omega(x)}{(x-x_i)\omega'(x_i)} \right]^2 dx > 0. \quad (5.5.4)$$

Второе замечание о коэффициентах A_k касается способов их вычисления. Для A_k было дано два явных выражения: (5.5.2) и (5.5.4). Можно указать иное выражение для A_k , более удобное для вычислений. Рассмотрим систему ортогональных на $\langle a, b \rangle$ по весу $p(x)$ многочленов

$$P_n(x) = a_n x^n + b_n x^{n-1} + \dots \quad (n=0, 1, \dots).$$

Для определенности формул, положим их нормированными:

$$\int_a^b p P_n^2 dx = 1, \quad a_n > 0.$$

Многочлен $P_n(x)$ отличается от $\omega(x)$ лишь численным множителем $P_n(x) = a_n \omega(x)$. Корни x_k многочлена $P_n(x)$ являются узлами квадратурного правила, коэффициенты же A_k имеют следующее выражение через P_n :

$$A_k = \int_a^b p(x) \frac{P_n(x)}{(x-x_k) P_n'(x_k)} dx.$$

Интеграл может быть просто вычислен, если воспользоваться известным в теории ортогональных многочленов тождеством Дарбу — Кристоффеля.*)

В нужной нам форме это соотношение может быть записано так:

$$(x-t) \sum_{i=0}^{n-1} P_i(x) P_i(t) = \frac{a_{n-1}}{a_n} [P_n(x) P_{n-1}(t) + P_{n-1}(x) P_n(t)].$$

Положим здесь $t = x_k$ и разделим обе части равенства на $x - x_k$:

$$\sum_{i=0}^{n-1} P_i(x) P_i(x_k) = \frac{a_{n-1}}{a_n} P_{n-1}(x_k) \frac{P_n(x)}{x - x_k}.$$

Умножим теперь обе части равенства на вес $p(x)$ и проинтегрируем по $\langle a, b \rangle$. Ввиду ортогональности между собой многочленов $P_k(x)$ и нормированности их интеграл

$$P_i(x_k) \int_a^b p(x) P_i(x) dx$$

*) См., например, Л. В. Гончаров. Теория интерполирования и приближения функций, гл. III, № 40. М., 1954.

равен нулю при $i \geq 1$ и равен единице при $i=0$.

$$1 = \frac{a_{n-1}}{a_n} P_{n-1}(x_k) \int_a^b p(x) \frac{P_n(x)}{x-x_k} dx = \frac{a_{n-1}}{a_n} P_n'(x_k) P_{n-1}(x_k) A_k, \quad (5.5.5)$$

$$A_k = \frac{a_n}{a_{n-1}} \cdot \frac{1}{P_n'(x_k) P_{n-1}(x_k)}.$$

Найденное выражение для A_k более удобно при вычислениях, так как не требует интегрирований.

5.5.3. Остаток квадратурного правила

Теорема 6. Если $p(x)$ сохраняет знак на $\langle a, b \rangle$ и f имеет непрерывную производную порядка $2n$ на $\langle a, b \rangle$, то существует такая точка $\xi \in \langle a, b \rangle$, что для остатка

$$R_n(f) = \int_a^b p(x) f(x) dx - \sum_{k=1}^n A_k f(x_k)$$

квадратурного правила наивысшей степени точности верно равенство

$$R_n(f) = \frac{1}{(2n)!} f^{(2n)}(\xi) \int_a^b p(x) \omega^2(x) dx. \quad (5.5.6)$$

Доказательство. Рассмотрим интерполяционный многочлен $H(x)$ степени не выше $2n-1$, удовлетворяющий условиям

$$H(x_k) = f(x_k), \quad H'(x_k) = f'(x_k).$$

При сделанном предположении о непрерывности $f^{(2n)}$, остаток интерполирования может быть представлен в форме

$$r(x) = \frac{1}{(2n)!} f^{(2n)}(\eta) \omega^2(x),$$

где η — некоторая точка отрезка, содержащего x и x_k ($k=1, \dots, n$), и

$$\int_a^b p(x) f(x) dx = \int_a^b p(x) H(x) dx + \frac{1}{(2n)!} \int_a^b f^{(2n)}(\eta) \omega^2(x) dx,$$

Существование последнего интеграла следует из существования двух других. Так как квадратурное правило верно для всех многочленов степени не выше $2n-1$, то

$$\int_a^b p(x)H(x)dx = \sum_{k=1}^n A_k H(x_k) = \sum_{k=1}^n A_k f(x_k)$$

и остаток $R_n(f)$ имеет, следовательно, значение

$$R_n(f) = \frac{1}{(2n)!} \int_a^b f^{(2n)}(\eta) p(x) \omega^2(x) dx.$$

Путем обычных рассуждений отсюда можно легко прийти к заключению о существовании точки ξ на $\langle a, b \rangle$, для которой выполняется равенство (5.5.6).

5.5.4. Сходимость квадратурного процесса наивысшей степени точности

Пусть весовая функция неотрицательна: $p(x) \geq 0$. Квадратурное правило наивысшей степени точности может быть построено для любого $n=1, 2, \dots$. Узлы и коэффициенты правила будут иметь свои значения для каждого n , и их мы будем обозначать x_k^n и A_k^n .

$$\int_a^b p(x)f(x)dx = \sum_{k=1}^n A_k^n f(x_k^n) + R_n(f) = Q_n(f) + R_n(f).$$

Говорят, что квадратурный процесс сходится для f , если

$$Q_n(f) \rightarrow \int_a^b p(x)f(x)dx \quad (n \rightarrow \infty).$$

Нашей задачей будет выяснить, для какого класса интегралов можно гарантировать сходимость квадратурного процесса наивысшей степени точности. Теорема, которую мы докажем сейчас, является частным случаем более общей теоремы, доказываемой ниже — в параграфе о сходимости квадратурных процессов. Но доказательство, которое мы приведем сейчас, основано на простых и хорошо известных фактах математического анализа и элементарно по ходу рассуждений, тогда как доказательство более общей теоремы опирается на значительно более сложные сведения из теории операторов.

Теорема 7. Если $p(x) \geq 0$, отрезок $[a, b]$ конечный и замкнутый и функция f непрерывна на нем, то квадратурный процесс наивысшей степени точности сходится.

Доказательство. Ввиду непрерывности f , при всяком $\varepsilon > 0$ существует многочлен $P(x)$ такой, что при любом $x \in [a, b]$ будет $|f(x) - P(x)| < \varepsilon$.

$$\begin{aligned} \left| \int_a^b p f dx - \sum_{k=1}^n A_k^n f(x_k^n) \right| &\leq \left| \int_a^b p f dx - \int_a^b p P dx \right| + \\ &+ \left| \int_a^b p P dx - \sum_{k=1}^n A_k^n P(x_k^n) \right| + \left| \sum_{k=1}^n A_k^n [P(x_k^n) - f(x_k^n)] \right|. \end{aligned}$$

Но

$$\left| \int_a^b p f dx - \int_a^b p P dx \right| = \left| \int_a^b p (f - P) dx \right| < \varepsilon \int_a^b p dx.$$

Кроме того, так как

$$\int_a^b p \cdot 1 dx = \sum_{k=1}^n A_k^n,$$

то

$$\left| \sum_{k=1}^n A_k^n [P(x_k^n) - f(x_k^n)] \right| < \varepsilon \sum_{k=1}^n A_k^n = \varepsilon \int_a^b p dx.$$

Наконец, если m есть степень многочлена P , то при $2n-1 \geq m$

$$\int_a^b p P dx = \sum_{k=1}^n A_k^n P(x_k^n)$$

и для таких n

$$\left| \int_a^b p f dx - \sum_{k=1}^n A_k^n f(x_k^n) \right| < 2\varepsilon \int_a^b p dx,$$

что доказывает теорему.

5.5.5. Замечание об интегрировании периодических функций

Мы закончим настоящий параграф замечанием об интегрировании гладких периодических функций. Оно отчасти выходит за границы параграфа, так как здесь речь будет идти о наивысшей тригонометрической, а не алгебраической степени точности.

Пусть $f(x)$ есть произвольная периодическая функция. Ее период всегда можно считать приведенным к 2π . Рассмотрим интеграл $\int_0^{2\pi} f(x) dx$, в котором весовая функция считается величиной постоянной, и для его вычисления будем строить правило вида

$$\int_0^{2\pi} f(x) dx \approx \sum_{k=1}^n A_k f(x_k), \quad 0 \leq x_k < 2\pi \quad (k=1, \dots, n). \quad (5.5.7)$$

В курсах анализа доказывается, что всякую непрерывную 2π -периодическую функцию f можно равномерно и сколь угодно точно приблизить при помощи тригонометрического многочлена

$$T_m(x) = a_0 + \sum_{k=1}^m (a_k \cos kx + b_k \sin kx).$$

Поэтому естественно стремиться параметры A_k и x_k выбрать так, чтобы правило (5.5.7) давало точный результат для многочленов $T_m(x)$ возможно высокой степени.

Можно просто проверить, что при любых A_k и x_k правило не может быть точным для всех тригонометрических многочленов степени n . Чтобы показать это, возьмем функцию

$$T(x) = \prod_{k=1}^n \sin^2 \frac{1}{2} (x - x_k).$$

Из равенства

$$\sin^2 \frac{1}{2} (x - x_k) = \frac{1}{2} [1 - \cos(x - x_k)]$$

ясно, что $T(x)$ есть тригонометрический многочлен степени n . Но для него правило (5.5.7) не может быть точным, так как $\int_0^{2\pi} T(x) dx > 0$, а $\sum_{k=1}^n A_k T(x_k) = 0$, ввиду того что все x_k являются корнями многочлена $T(x)$.

Тригонометрическая степень точности правила (5.5.7) всегда меньше n и при помощи выбора A_k и x_k ее можно надеяться сделать равной самое большее $n-1$. Как оказывается, наивысшая степень точности $n-1$ достигается квадратурной формулой с равными коэффициентами $A_k = \frac{2\pi}{n}$ ($k=1, \dots, n$) и равноотстоящими узлами.

Пусть α есть любое число, выполняющее неравенство $0 \leq \alpha < h = \frac{2\pi}{n}$. Рассмотрим точки $x_i = \alpha + ih$ ($i=0, 1, \dots, n-1$). Они лежат на отрезке $0 \leq x < 2\pi$. Примем их за узлы x_k и построим квадратурное правило

$$\int_0^{2\pi} f(x) dx \approx \frac{2\pi}{n} \sum_{k=1}^n f\left[\alpha + (k-1) \frac{2\pi}{n}\right]. \quad (5.5.8)$$

Убедимся в том, что оно является точным для всех тригонометрических многочленов степени $n-1$. Для этого достаточно проверить, что (5.5.8) точно выполняется для функций e^{imx} ($m=0, 1, \dots, n-1$).

При $m=0$ это, очевидно, верно. Для $1 \leq m \leq n-1$ вычисления дают результаты:

$$\begin{aligned} \int_0^{2\pi} e^{imx} dx &= \frac{1}{im} (e^{im 2\pi} - 1) = 0, \\ \sum_{k=1}^m e^{im[\alpha + (k-1)h]} &= e^{im\alpha} \sum_{k=1}^m e^{i(k-1)mh} = \\ &= e^{im\alpha} \frac{e^{imnh} - 1}{e^{imh} - 1} = e^{im\alpha} \frac{e^{im 2\pi} - 1}{e^{imh} - 1} = 0, \end{aligned}$$

что доказывает точное выполнение (5.5.8) и в этом случае.

§ 5.6. НЕКОТОРЫЕ ЧАСТНЫЕ СЛУЧАИ КВАДРАТУРНЫХ ПРАВИЛ НАИВЫСШЕЙ АЛГЕБРАИЧЕСКОЙ СТЕПЕНИ ТОЧНОСТИ

Ниже будут рассматриваться квадратурные правила, отвечающие весовым функциям, особенно часто встречающимся в приложениях.

5.6.1. Постоянная весовая функция

Отрезок интегрирования $[a, b]$ считается конечным, и интеграл берется в форме

$$\int_a^b f(x) dx,$$

где $f(x)$ предполагается достаточно гладкой функцией. Интегралы такого типа особенно часто встречаются в приложениях. Соответствующее этому случаю правило было найдено Гауссом и носит его имя.

Всякий конечный отрезок $[a, b]$ линейной заменой переменной может быть преобразован в $[-1, 1]$, и мы будем считать, что интеграл приведен к виду

$$\int_{-1}^1 f(x) dx. \quad (5.6.1)$$

Систему многочленов, ортогональную на $[-1, 1]$ с весом $\rho(x) \equiv 1$, образуют, как известно, многочлены Лежандра

$$P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n.$$

В квадратурной формуле с n узлами

$$\int_{-1}^1 f(x) dx \approx \sum_{k=1}^n A_k f(x_k), \quad (5.6.2)$$

имеющей наивысшую степень точности $2n-1$, узлы x_k ($k=1, 2, \dots, n$) должны располагаться в корнях многочлена Лежандра степени n :

$$P_n(x_k) = 0 \quad (k=1, 2, \dots, n).$$

Коэффициенты A_k могут быть вычислены, например, при помощи равенства вида (5.5.5). Напомним, что при записи этого равенства мы пользовались нормированными многочленами. Поэтому при применении его в нашем случае мы должны воспользоваться многочленами

$$\rho_n(x) = \sqrt{\frac{2n+1}{2}} P_n(x).$$

Старшие коэффициенты их имеют значение

$$a_n = \sqrt{\frac{2n+1}{2}} \frac{(2n)!}{2^n (n!)^2}.$$

Несложные вычисления позволяют найти для A_k равенство

$$A_k = \frac{2}{nP_{n-1}(x_k)P_n'(x_k)}. \quad (5.6.3)$$

Оно может быть приведено к виду, несколько более удобному для вычислений, если воспользоваться известным в теории многочленов Лежандра соотношением

$$(1-x^2)P_n'(x) = n[P_{n-1}(x) - xP_n(x)].$$

Положим здесь $x = x_k$:

$$(1-x_k^2)P_n'(x_k) = nP_{n-1}(x_k).$$

Это равенство позволяет исключить одну из величин P_n' или P_{n-1} и привести (5.6.3), например, к виду

$$A_k = \frac{2}{n^2} \cdot \frac{1-x_k^2}{P_{n-1}^2(x_k)}. \quad (5.6.4)$$

Когда функция f имеет на $[-1, 1]$ непрерывную производную порядка $2n$, для нахождения погрешности гауссовой формулы (5.6.2) можно воспользоваться равенством (5.5.6). В нем мы должны положить $p(x) \equiv 1$. Что же касается многочлена $\omega(x)$, то он может отличаться от $P_n(x)$ только постоянным множителем и, так как старший коэффициент $\omega(x)$ равен единице, а в $P_n(x)$ он равен

$$\frac{(2n)!}{2^n (n!)^2},$$

то

$$\omega(x) = \frac{2^n (n!)^2}{(2n)!} P_n(x).$$

Кроме того, если принять во внимание, что

$$\int_{-1}^1 P_n^2 dx = \frac{2}{2n+1},$$

из (5.5.6) для погрешности правила Гаусса получится

$$R_n(f) = \frac{2^{2n+1}}{(2n+1)(2n)!} \left[\frac{(n!)^2}{(2n)!} \right]^2 f^{(2n)}(\xi) \quad (-1 \leq \xi \leq 1). \quad (5.6.5)$$

5.6.2. Интегралы вида $\int_a^b (b-x)^\alpha (x-a)^\beta f(x) dx$

Как видно из записи интеграла, соответствующее ему правило предназначено для интегрирования функций, имеющих на концах отрезка $[a, b]$ степенные особенности, или когда заранее известно, что точки a и b являются нулями f , и известна также кратность этих нулей.

Линейным преобразованием $x = \frac{1}{2}(a+b) + \frac{1}{2}(b-a)t$ отрезок $[a, b]$ приводится к $[-1, 1]$, и можно ограничиться рассмотрением интеграла

$$\int_{-1}^1 (1-x)^\alpha (1+x)^\beta f(x) dx \quad (\alpha, \beta > -1).$$

Системой многочленов, ортогональных на $[-1, 1]$ по весу

$$\rho(x) = (1-x)^\alpha (1+x)^\beta,$$

является система многочленов Якоби

$$P_n^{(\alpha, \beta)}(x) = \frac{(-1)^n}{2^n n!} (1-x)^{-\alpha} (1+x)^{-\beta} \times \\ \times \frac{d^n}{dx^n} [(1-x)^{\alpha+n} (1+x)^{\beta+n}] = \frac{\Gamma(\alpha+\beta+2n+1)}{2^n n! \Gamma(\alpha+\beta+n+1)} x^n + \dots \quad (5.6.6)$$

При построении для рассматриваемого интеграла правила наивысшей степени точности $2n-1$

$$\int_{-1}^1 (1-x)^\alpha (1+x)^\beta f(x) dx \approx \sum_{k=1}^n A_k f(x_k) \quad (5.6.7)$$

мы должны в качестве узлов x_k взять корни многочлена Якоби степени n :

$$P_n^{(\alpha, \beta)}(x_k) = 0 \quad (k=1, 2, \dots, n).$$

Коэффициенты A_k могут быть найдены при помощи (5.5.5). Нормированные многочлены Якоби, как известно,^{*)} есть

$$P_n^{(\alpha, \beta)}(x) = \delta_n^{-\frac{1}{2}} P_n^{(\alpha, \beta)}(x),$$

^{*)} Г. Сегё. Ортогональные многочлены. М., 1962.

где

$$\delta_n = \frac{2^{\alpha+\beta+1} \Gamma(\alpha+n+1) \Gamma(\beta+n+1)}{(\alpha+\beta+2n+1)n! \Gamma(\alpha+\beta+n+1)}.$$

Старший коэффициент его

$$a_n = \delta_n \cdot \frac{1}{2} \cdot \frac{\Gamma(\alpha+\beta+2n+1)}{2^n n! \Gamma(\alpha+\beta+n+1)}.$$

Несложные вычисления дадут для A_k значение

$$A_k = \frac{(\alpha+\beta+2n) 2^{\alpha+\beta} \Gamma(\alpha+n) \Gamma(\beta+n)}{n! \Gamma(\alpha+\beta+n+1) P_{n-1}^{(\alpha, \beta)}(x_k) [P_n^{(\alpha, \beta)}(x_k)]'}.$$

Его можно упростить, если воспользоваться известным в теории многочленов Якоби равенством.*)

$$\begin{aligned} & (\alpha+\beta+2n)(1-x^2) \frac{d}{dx} P_n^{(\alpha, \beta)}(x) = \\ & = -n[(\alpha+\beta+2n)x+\beta-\alpha] P_n^{(\alpha, \beta)}(x) + 2(\alpha+n)(\beta+n) P_{n-1}^{(\alpha, \beta)}(x). \end{aligned}$$

Положив здесь $x=x_k$, получим равенство

$$(\alpha+\beta+2n)(1-x_k^2) [P_n^{(\alpha, \beta)}]' = 2(\alpha+n)(\beta+n) P_{n-1}^{(\alpha, \beta)}(x_k),$$

позволяющее исключить любую из величин $P_n^{(\alpha, \beta)'}(x_k)$ или $P_{n-1}^{(\alpha, \beta)}(x_k)$.

Например, если мы исключим $P_{n-1}^{(\alpha, \beta)}(x_k)$, то для A_k получится

$$A_k = 2^{\alpha+\beta+1} \frac{\Gamma(\alpha+n+1) \Gamma(\beta+n+1)}{n! \Gamma(\alpha+\beta+n+1) (1-x_k^2) [P_n^{(\alpha, \beta)'}(x_k)]^2}. \quad (5.6.7')$$

Для построения остатка правила (5.6.7) можно, как и в случае гауссова правила, воспользоваться равенством (5.5.6).

Старший коэффициент $P_n^{(\alpha, \beta)}(x)$ указан в (5.6.6). Многочлен $\omega(x)$ связан с $P_n^{(\alpha, \beta)}$ равенством

*) См. предыдущую сноску.

$$\omega(x) = \frac{2^n n! \Gamma(\alpha + \beta + n + 1)}{\Gamma(\alpha + \beta + 2n + 1)} P_n^{(\alpha, \beta)}(x).$$

Поэтому

$$\begin{aligned} R_n(f) &= \frac{f^{(2n)}(\xi)}{(2n)!} \left[\frac{2^n n! \Gamma(\alpha + \beta + n + 1)}{\Gamma(\alpha + \beta + 2n + 1)} \right]^2 \times \\ &\quad \times \int_{-1}^1 (1-x)^\alpha (1+x)^\beta [P_n^{(\alpha, \beta)}(x)]^2 dx = \\ &= \frac{f^{(2n)}(\xi)}{(2n)!} \frac{2^{\alpha+\beta+2n+1} n! \Gamma(\alpha + n + 1) \Gamma(\beta + n + 1) \Gamma(\alpha + \beta + n + 1)}{(\alpha + \beta + 2n + 1) \Gamma^2(\alpha + \beta + 2n + 1)} \\ &\quad (-1 \leq \xi \leq 1). \end{aligned}$$

Квадратурная формула (5.6.7) содержит произвольные параметры α и β и является источником многих полезных частных случаев. Частным случаем ее является квадратурное правило Гаусса, получающееся при $\alpha = \beta = 0$. Рассмотрим еще частный случай, когда $\alpha = \beta = -\frac{1}{2}$. Весовая функция здесь будет

$$\rho(x) = \frac{1}{\sqrt{1-x^2}}.$$

Многочлены Якоби $P_n^{(-0,5; -0,5)}(x)$ только численным множителем отличаются от многочленов Чебышева первого рода:

$$P_n^{(-0,5; -0,5)}(x) = C_n T_n(x) = C_n \cos(n \arccos x).$$

Узлы квадратурного правила должны совпадать с нулями многочлена T_n :

$$x_k = \cos \frac{2k-1}{2n} \pi \quad (k=1, 2, \dots, n).$$

Коэффициенты A_k можно подсчитать при помощи (5.6.7'):

$$T_n'(x_k) = \sin(n \arccos x_k) \frac{n}{\sqrt{1-x_k^2}} = \frac{(-1)^{k-1} n}{\sqrt{1-x_k^2}},$$

$$(1-x_k^2) [P_n^{(-0,5; -0,5)'}(x_k)]^2 = C_n^2 (1-x_k^2) T_n'^2(x_k) = C_n^2 n^2,$$

$$A_k = \frac{\Gamma^2\left(n + \frac{1}{2}\right)}{n^2 \Gamma(n) C_n^2 n^2}.$$

Правая часть последнего равенства не зависит от k , и все коэффициенты A_k будут одинаковы. Общую величину их обозначим A . Наиболее просто A можно найти, если воспользоваться тем, что квадратурная формула должна быть точной для $f \equiv 1$:

$$\sum_{k=1}^n A_k = nA = \int_{-1}^1 \frac{dx}{\sqrt{1-x^2}} = \pi, \quad A = \frac{\pi}{n}.$$

Таким образом, квадратурное правило наивысшей степени точности с весом $p(x) = \frac{1}{\sqrt{1-x^2}}$ имеет вид

$$\int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}} dx = \frac{\pi}{n} \sum_{k=1}^n f\left(\cos \frac{2k-1}{2n} \pi\right) + R_n(f), \quad (5.6.8)$$

$$R_n(f) = \frac{\pi}{2^{2n-1} (2n)!} f^{(2n)}(\xi) \quad (-1 \leq \xi \leq 1).$$

В связи с этим П. Л. Чебышевым была поставлена задача о построении правил приближенной квадратуры с равными коэффициентами с любой весовой функцией $p(x)$, в частности с постоянным весом. Эта задача будет рассматриваться в одном из следующих параграфов.

Укажем еще на возможность приложения правила квадратур с весом Якоби (5.6.7) к задаче кратного численного интегрирования.

Проблема вычисления кратных интегралов существенно отличается от случая простых интегралов. Если в однократном интеграле практически важная область интегрирования очень простая — ею является отрезок, то для многократных интегралов область интегрирования может быть очень сложной. Это обстоятельство сильно затрудняет задачу построения правил интегрирования, и в одном из способов ее решения интегралы с весовой функцией Якоби могут оказать, как будет показано ниже, существенную помощь. Идея этого с достаточной полнотой может быть выяснена на примере двойного интеграла в декартовых координатах

$$I = \iint_{\omega} f(x, y) dx dy.$$

Если область интегрирования ω обладает хорошо известными в курсах анализа геометрическими свойствами, вычисление двойного интеграла может быть приведено к нахождению двух однократных интегралов

$$F(x) = \int_{y_1(x)}^{y_2(x)} f(x, y) dy, \quad I = \int_a^b F(x) dx,$$

где $y_1(x)$, $y_2(x)$, a , b зависят от формы области и имеют известные значения.

Выбор правила для вычисления I должен быть согласован со свойствами функции f , во-первых, и со свойствами области ω , во-вторых. Сейчас мы хотим выяснить вопрос о влиянии на выбор правила свойств ω . Поэтому будем предполагать f достаточно гладкой всюду в ω .

В интеграле $F(x)$ функция $f(x, y)$, по нашему допущению, не имеет особенностей, и он может быть вычислен по одному из известных правил с постоянным весом для простых интегралов, например по правилу Гаусса, Симпсона и др. Форма области ω оказывает влияние только на границы интегрирования y_1 и y_2 . Отрезок $[y_1, y_2]$ можно привести к каноническому, например $[0, 1]$, подстановкой $y = y_1 + (y_2 - y_1)\eta$ ($0 \leq \eta \leq 1$). Тогда получим

$$\begin{aligned} F(x) &= [y_2(x) - y_1(x)] \int_0^1 f\{x, y_1(x) + [y_2(x) - y_1(x)]\eta\} d\eta = \\ &= [y_2(x) - y_1(x)] \Phi(x). \end{aligned}$$

Выделившийся при замене в интеграле

$$I = \int_a^b F(x) dx = \int_a^b [y_2(x) - y_1(x)] \Phi(x) dx \quad (5.6.9)$$

множитель $y_2(x) - y_1(x)$ является естественной весовой функцией в (5.6.9). При вычислении (5.6.9) можно воспользоваться любым квадратурным правилом, построенным для веса $p(x) = y_2(x) - y_1(x)$, например правилом наивысшей степени точности.

Такой полный учет формы области, вероятно, не разумно делать, так как каждой области ω будет отвечать свой вес $p(x)$ и пришлось бы вычислять большое число таблиц квадратурных узлов x_k и коэффициентов A_k . Можно упростить задачу на основании следующих простых соображений. Рассмотрим две весовые функции $p(x)$ и $q(x)$, отличающиеся друг от друга достаточно гладким множителем $\rho(x)$, не обращающимся в нуль на $[a, b]$: $q(x) = \rho(x)p(x)$.

Можно ожидать, что квадратурные правила, соответствующие этим двум весам $p(x)$ и $q(x)$, будут близкими по своей точности.

Возвратимся к весовой функции Якоби $q(x) = (x-a)^\beta (b-x)^\alpha$. Она зависит от двух показателей α и β , и их часто можно подобрать так, чтобы отношение

$$\rho(x) = \frac{y_2(x) - y_1(x)}{(b-x)^\alpha (x-a)^\beta} \quad (a \leq x \leq b)$$

было ограничено сверху и снизу положительными числами $0 < m \leq \rho(x) \leq M < \infty$. В этом случае для вычисления интеграла (5.6.9) можно воспользоваться весом $q(x) = (b-x)^\alpha \times (x-a)^\beta$, преобразовав интеграл I к виду

$$I = \int_a^b (b-x)^\alpha (x-a)^\beta \Psi(x) dx, \quad \Psi(x) = \rho(x) \Phi(x),$$

и известными таблицами x_k и A_k для яковиева веса.

Приведем пояснительный пример. Допустим, что область ω имеет вид, изображенный на рис. 5.6.1, и ее контур λ в точках A и B имеет с прямыми $x=a$ и $x=b$ соприкосновение первого порядка.*) В качестве $q(x)$ можно тогда взять $q(x) = \sqrt{(b-x)(x-a)}$ и привести интеграл I к виду

$$I = \int_a^b \sqrt{(b-x)(x-a)} \Psi(x) dx, \quad \Psi(x) = [(b-x)(x-a)]^{-\frac{1}{2}} F(x).$$

5.6.3. Интегралы вида $\int_0^\infty x^\alpha e^{-x} f(x) dx$

Ортогональными на полуоси $[0, \infty)$ по весу $p(x) = x^\alpha e^{-x}$ ($\alpha > -1$) являются многочлены Чебышева — Лягерра

$$L_n^{(\alpha)}(x) = (-1)^n x^{-\alpha} e^x \frac{d^n}{dx^n} (x^{\alpha+n} e^{-x}) = x^n - \frac{n(n+\alpha)}{1!} x^{n-1} + \dots$$

В квадратурном правиле наивысшей степени точности

*) Говорят, что точка A является точкой соприкосновения первого порядка, если уравнение контура λ вблизи A можно записать в форме $x = a + c_2(y-y_0)^2 + c_3(y-y_0)^3 + \dots$, $c_2 \neq 0$. Аналогично для точки B .

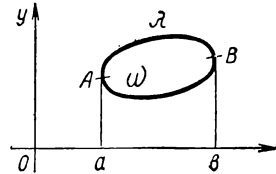


Рис. 5.6.1

$$\int_0^{\infty} x^{\alpha} e^{-x} f(x) dx = \sum_{k=1}^n A_k f(x_k) + R_n(f) \quad (5.6.10)$$

узлами x_k должны служить корни многочлена $L_n^{(\alpha)}(x)$ степени n :

$$L_n^{(\alpha)}(x_k) = 0 \quad (k=1, 2, \dots, n).$$

Для $L_n^{(\alpha)}(x)$ верно равенство

$$\int_0^{\infty} x^{\alpha} e^{-x} [L_n^{(\alpha)}(x)]^2 dx = n! \Gamma(n+\alpha+1),$$

поэтому ортонормальными многочленами Чебышева — Лягерра будут

$$l_n^{(\alpha)}(x) = [n! \Gamma(n+\alpha+1)]^{-\frac{1}{2}} L_n^{(\alpha)}(x).$$

Старшие коэффициенты их

$$a_n = [n! \Gamma(n+\alpha+1)]^{-\frac{1}{2}}$$

и формула (5.5.5) для A_k дает значение

$$A_k = \frac{\Gamma(n) \Gamma(n+\alpha)}{L_n^{(\alpha)'}(x_k) L_{n-1}^{(\alpha)}(x_k)}.$$

Для многочленов Лягерра известно следующее соотношение:

$$x L_n^{(\alpha)'}(x) = n L_n^{(\alpha)}(x) + n(n+\alpha) L_{n-1}^{(\alpha)}(x).$$

Если в нем положить $x = x_k$, то получится равенство

$$x_k L_n^{(\alpha)'}(x_k) = n(n+\alpha) L_{n-1}^{(\alpha)}(x_k),$$

позволяющее найденное выражение для A_k привести к виду

$$A_k = \frac{\Gamma(n+1) \Gamma(n+\alpha+1)}{x_k [L_n^{(\alpha)'}(x_k)]^2}.$$

Когда f имеет непрерывную производную порядка $2n$ на полуоси $[0, \infty)$, равенство (5.5.6) дает для остатка $R_n(f)$ в (5.6.10) следующее представление:

$$R_n(f) = \frac{\Gamma(n+1)\Gamma(\alpha+n+1)}{(2n)!} f^{(2n)}(\xi), \quad \xi \in [0, \infty). \quad (5.6.11)$$

5.6.4. Интегралы вида $\int_{-\infty}^{\infty} e^{-x^2} f(x) dx$

Систему многочленов, ортогональных на оси $-\infty < x < \infty$ по весу e^{-x^2} , образуют многочлены Чебышева — Эрмита

$$H_n(x) = (-1)^n e^{x^2} \frac{d^n}{dx^n} e^{-x^2} = 2^n x^n - \dots$$

В соответствующей квадратурной формуле наивысшей алгебраической степени точности $2n-1$

$$\int_{-\infty}^{\infty} e^{-x^2} f(x) dx = \sum_{k=1}^n A_k f(x_k) + R_n(f) \quad (5.6.12)$$

узлы x_k должны быть корнями многочлена H_n :

$$H_n(x_k) = 0 \quad (k=1, 2, \dots, n).$$

Так как

$$\int_{-\infty}^{\infty} e^{-x^2} H_n^2(x) dx = 2^n n! \sqrt{\pi},$$

нормированными многочленами Чебышева — Эрмита являются

$$h_n(x) = [2^n n! \sqrt{\pi}]^{-\frac{1}{2}} H_n(x).$$

Старшие коэффициенты их есть

$$a_n = 2^{\frac{n}{2}} [n! \sqrt{\pi}]^{-\frac{1}{2}}.$$

Приняв, кроме того, во внимание соотношение $H_n'(x) = 2nH_{n-1}(x)$, при помощи (5.5.5) для A_k получим

$$A_k = \frac{2^{n+1}n! \sqrt{\pi}}{H_n^2(x_k)}. \quad (5.6.13)$$

Наконец, при предположении существования у f непрерывной производной порядка $2n$ на всей оси $-\infty < x < \infty$ для остатка $R_n(f)$, на основании (5.5.6), получим

$$R_n(f) = \frac{n! \sqrt{\pi}}{2^n (2n)!} f^{(2n)}(\xi). \quad (5.6.14)$$

§ 5.7. КВАДРАТУРНЫЕ ПРАВИЛА НАИВЫСШЕЙ СТЕПЕНИ ТОЧНОСТИ, ИМЕЮЩИЕ ФИКСИРОВАННЫЕ ЗАРАНЕЕ УЗЛЫ

5.7.1. Некоторые общие теоремы

При вычислениях нередко приходится иметь дело с интегралами, в которых заранее известны или легко вычисляются значения интегрируемой функции в одной или нескольких точках. Примером может служить интеграл, содержащий решение граничной задачи, когда значения функции на концах отрезка задаются заранее. Поэтому естественно строить квадратурные правила, которые позволяли бы учитывать эти известные значения, если не всегда, то, по крайней мере, в наиболее часто встречающихся и важных случаях.

Рассмотрим правило квадратур вида

$$\int_a^b p(x)f(x)dx \approx \sum_{k=1}^n A_k f(x_k) + \sum_{i=1}^m B_i f(a_i), \quad (5.7.1)$$

содержащее m фиксированных узлов a_1, a_2, \dots, a_m . Формула содержит $2n+m$ параметров A_k, x_k ($k=1, \dots, n$) и B_i ($i=1, \dots, m$). Выясним возможность такого выбора их, чтобы сделать равенство (5.7.1) точным для всевозможных алгебраических многочленов степени $2n+m-1$.

Напомним, что при любых x_k и a_i только за счет выбора коэффициентов A_k и B_i равенство (5.7.1) можно сделать точным для всяких многочленов степени $n+m-1$. Для этого достаточно считать правило интерполяционным.

В нашем случае это означает, что его коэффициенты должны быть следующими:

$$A_k = \int_a^b p(x) \frac{\omega(x)\Omega(x)}{(x-x_k)\omega'(x_k)\Omega(x_k)} dx,$$

$$B_i = \int_a^b p(x) \frac{\omega(x)\Omega(x)}{(x-a_i)\omega(a_i)\Omega'(a_i)} dx, \quad (5.7.2)$$

$$\omega(x) = (x-x_1)\dots(x-x_n), \quad \Omega(x) = (x-a_1)\dots(x-a_m).$$

После этого в нашем распоряжении останется еще выбор узлов x_k ($k=1, \dots, n$).

Теорема 1. Для того чтобы правило (5.7.1) было точным для многочленов степени $2n+m-1$, необходимо и достаточно выполнение двух условий:

1) правило является интерполяционным, т. е. его коэффициенты A_k и B_i имеют значения (5.7.2);

2) многочлен $\omega(x)$ ортогонален на $[a, b]$ по весу $p(x)\Omega(x)$ ко всякому многочлену $Q(x)$ степени меньшей n :

$$\int_a^b p(x)\Omega(x)\omega(x)Q(x)dx=0. \quad (5.7.3)$$

Доказательство. Необходимость первого условия легко проверить, так как если правило точно для многочленов степени $2n+m-1$, то оно точно и для многочленов степени $n+m-1$, а тогда по теореме 1 § 5.2, оно должно быть интерполяционным. Для доказательства необходимости второго условия достаточно положить $f=\omega(x)\Omega(x)Q(x)$. Так как f есть многочлен степени не выше $2n+m-1$, для него равенство (5.7.1) должно выполняться точно. Но, ввиду того что f обращается в нуль в точках x_k и a_i , правая часть (5.7.1) равна нулю и должно выполняться (5.7.3).

Пусть f есть произвольный многочлен степени $2n+m-1$. Если разделить f на $\omega\Omega$, его можно представить в форме

$$f(x) = \Omega(x)\omega(x)Q(x) + r(x),$$

где $Q(x)$ и $r(x)$ — многочлены степеней соответственно не больше $n-1$ и $n+m-1$. Очевидно, $f(x_k) = r(x_k)$ и $f(a_i) = r(a_i)$.

Если выполняются условия 1 и 2, то будет верной следующая цепь точных равенств, доказывающая достаточность условий теоремы:

$$\begin{aligned} \int_a^b p f dx &= \int_a^b p \Omega \omega Q dx + \int_a^b p r dx = \int_a^b p r dx = \\ &= \sum_{k=1}^n A_k r(x_k) + \sum_{i=1}^m B_i r(a_i) = \sum_{k=1}^n A_k f(x_k) + \sum_{i=1}^m B_i f(a_i). \end{aligned}$$

Построение правила (5.7.1), точного для многочленов степени $2n+m-1$, приводится к нахождению многочлена $\omega(x)$, удовлетворяющего условию ортогональности (5.7.3). Корни его x_k должны быть приняты за узлы x_k правила. Желательно, чтобы корни принадлежали отрезку интегрирования $[a, b]$, так как правила квадратур с узлами, лежащими вне отрезка интегрирования, имеют сравнительно ограниченную область применения.

Допустим, что многочлен $\omega(x)$, имеющий нужные нам свойства, существует и правило (5.7.1), точное для многочленов степени $2n+m-1$, может быть построено. Получим одно из возможных представлений погрешности правила. Выполним интерполирование f при помощи многочлена $H(x)$ степени $2n+m-1$ по условиям

$$H(a_i) = f(a_i) \quad (i=1, 2, \dots, m), \quad H(x_k) = f(x_k), \\ H'(x_k) = f'(x_k) \quad (k=1, \dots, n).$$

При предположении о существовании у f непрерывной производной порядка $2n+m$ остаток интерполирования $r(x) = f(x) - H(x)$ можно записать в виде

$$r(x) = \omega^2(x) \Omega(x) \frac{f^{(2n+m)}(\xi)}{(2n+m)!} \quad (a < \xi < b).$$

Для погрешности $R(f)$ квадратуры верно равенство $R(f) = R(H) + R(r)$ и, так как $R(H) = 0$, будет

$$R(f) = R(r) = \int_a^b p r dx - \sum_{k=1}^n A_k r(x_k) - \sum_{i=1}^m B_i r(a_i).$$

Ввиду $r(x_k) = 0$ ($k=1, \dots, n$) и $r(a_i) = 0$ ($i=1, \dots, m$) для $R(f)$ получим

$$R(f) = \int_a^b p(x) r(x) dx = \frac{1}{(2n+m)!} \int_a^b p(x) \Omega(x) \omega^2(x) f^{(2n+m)}(\xi) dx. \quad (5.7.4)$$

Полученное выражение для $R(f)$ позволяет просто решить вопрос о степени точности (5.7.1). Покажем, что если

$$I = \int_a^b p \Omega \omega^2 dx \neq 0,$$

то равенство (5.7.1) не может быть точным для многочленов степени $2n+m$ и, стало быть, степень точности его равна $2n+m-1$. В самом деле, когда f есть многочлен степени $2n+m$, производная $f^{(2n+m)}$ будет величиной постоянной, не равной нулю. Для такого многочлена остаток $R(f)$ имеет значение

$$R(f) = \frac{f^{(2n+m)}}{(2n+m)!} \int_a^b p \Omega \omega^2 dx,$$

отличное от нуля, и равенство (5.7.1) не может быть точным.

В исключительном случае, когда $l=0$, на рассмотрении которого мы не останавливаемся, степень точности правила больше $2n+m-1$ и может быть указан признак для ее определения.

5.7.2. Некоторые частные квадратурные правила

Рассмотрим частные случаи правил интегрирования с постоянной весовой функцией $p(x) \equiv 1$ и одним или двумя фиксированными узлами, лежащими на концах отрезка интегрирования. Последний считается конечным. Чтобы воспользоваться при вычислениях многочленами Якоби, будем считать этот отрезок приведенным к $[-1, 1]$. Полагая $m=1$, возьмем квадратурную формулу вида

$$\int_{-1}^1 f(x) dx = A f(-1) + \sum_{k=1}^n A_k f(x_k) + R(f). \quad (5.7.5)$$

Здесь $\Omega(x) = 1+x$. Вспомогательная весовая функция, участвующая в условии ортогональности (5.7.3): $\rho(x) = p(x)\Omega(x) = 1+x$ положительна внутри $[-1, 1]$. Многочлен $\omega(x)$ существует при всяких значениях $n=1, 2, \dots$. Он ортогонален на $[-1, 1]$ по весу $\rho(x) = 1+x$ ко всякому многочлену меньшей степени и может, следовательно, отличаться от многочлена Якоби $P_n^{(0,1)}(x)$ лишь численным множителем, равным обратной величине старшего коэффициента,

$$\omega(x) = \frac{2^n n! \Gamma(n+2)}{\Gamma(2n+2)} P_n^{(0,1)}(x).$$

Наивысшая степень точности формулы (5.7.5) равна $2n$. Она достигается, если в качестве узлов x_k взять корни многочлена $P_n^{(0,1)}(x)$ и коэффициенты определить при помощи равенств (5.7.2), которые в рассматриваемом случае принимают вид:

$$A_k = \frac{1}{1+x_{k-1}} \int_{-1}^1 (1+x) \frac{\omega(x)}{(x-x_k)\omega'(x_k)} dx \quad (k=1, \dots, n),$$

$$A = [P_n^{(0,1)}(-1)]^{-1} \int_{-1}^1 P_n^{(0,1)}(x) dx.$$

Оба интеграла могут быть без больших затруднений вычислены при помощи известных фактов теории многочленов Якоби, на которых мы не останавливаемся, и окончательно получится

$$A_k = \frac{4}{(1+x_k)(1-x_k^2)[P_n^{(0,1)}(x_k)]^2}, \quad (5.7.6)$$

$$A = \frac{2}{(n+1)^2}.$$

Остаток $R(f)$ может быть найден при помощи общей формулы (5.7.4):

$$R(f) = \frac{1}{(2n+1)!} \int_{-1}^1 (1+x)\omega^2(x)f^{(2n+1)}(\xi)dx.$$

Так как ядро интеграла $(1+x)\omega^2(x)$ сохраняет знак на $[-1, 1]$, на этом отрезке найдется такая точка η , что будет верным равенство

$$\begin{aligned} R(f) &= \frac{f^{(2n+1)}(\eta)}{(2n+1)!} \int_{-1}^1 (1+x)\omega^2(x)dx = \\ &= \frac{f^{(2n+1)}(\eta)}{(2n+1)!} \left[\frac{2^n n! (n+1)!}{(2n+1)!} \right]^2 \int_{-1}^1 (1+x)[P_n^{(0,1)}(x)]^2 dx = \\ &= \frac{2}{n+1} \left[\frac{2^n n! (n+1)!}{(2n+1)!} \right]^2 \frac{f^{(2n+1)}(\eta)}{(2n+1)!} \quad (-1 < \eta < 1). \end{aligned} \quad (5.7.7)$$

Случай фиксированного узла в точке 1 и соответствующей ему формулы

$$\int_{-1}^1 f(x)dx = \sum_{k=1}^n A_k f(x_k) + A f(1) + R(f)$$

приводится к (5.7.5) заменой x на $-x$ и рассматриваться отдельно не будет.

Остановимся еще на правиле интегрирования с двумя фиксированными узлами в точках -1 и 1 :

$$\int_{-1}^1 f(x) dx = Af(-1) + \sum_{k=1}^n A_k f(x_k) + Bf(1) + R(f). \quad (5.7.8)$$

Здесь $\Omega(x) = 1 - x^2$. Вспомогательная весовая функция

$$\rho(x) = p(x)\Omega(x) = 1 - x^2$$

положительна на $(-1, 1)$, и многочлен $\omega(x)$ существует при всяком n . Он ортогонален на $[-1, 1]$ по весу $1 - x^2$ ко всякому многочлену низшей степени и отличается от якобиева многочлена $P_n^{(1,1)}(x)$ постоянным множителем

$$\omega(x) = \frac{2^n n! \Gamma(n+3)}{\Gamma(2n+3)} P_n^{(1,1)}(x).$$

Наивысшая степень точности (5.7.8) равна $2n+1$. Она достигается, если за x_k принять корни $P_n^{(1,1)}(x)$ и коэффициенты A_k, A, B вычислить согласно (5.7.2). Расчеты дадут для коэффициентов и остатка значения:

$$\left. \begin{aligned} A_k &= 8 \cdot \frac{n+1}{n+2} \cdot \frac{1}{(1-x_k^2) [P_n^{(1,1)'}(x_k)]^2}, \quad A=B=\frac{2}{(n+1)(n+2)}, \\ R(f) &= \frac{8(n+1)}{(2n+3)(n+2)} \left[\frac{2^n n! (n+2)!}{(2n+2)!} \right]^2 \frac{f^{(2n+2)}(\eta)}{(2n+2)!} \\ &\quad (-1 < \eta < 1). \end{aligned} \right\} \quad (5.7.9)$$

§ 5.8. КВАДРАТУРНЫЕ ПРАВИЛА С РАВНЫМИ КОЭФФИЦИЕНТАМИ

5.8.1. Построение формул Чебышева. Существование и единственность

Правила приближенного интегрирования, все коэффициенты которых одинаковы:

$$\int_a^b p(x) f(x) dx \approx C_n \sum_{k=1}^n f(x_k), \quad (5.8.1)$$

весьма удобны при графических расчетах, так как сумму ординат легко можно снять с чертежа при помощи простых измерительных приборов, таких, например, как длинномер.

Теорема 1. Если $\int_a^b p(x)dx \neq 0$, формула (5.8.1) с действительными или комплексными узлами x_k всегда может быть построена и при этом единственным способом.

Когда среди x_k существуют комплексные, правило (5.8.1) будет иметь ограниченное значение и может оказаться полезным лишь при интегрировании аналитических функций f , регулярных в области, охватывающей отрезок $[a, b]$ и достаточно широкой.

Поэтому одной из задач теории формул Чебышева является нахождение тех случаев, когда все узлы x_k будут действительными.

Правило (5.6.8) интегрирования с весом $p(x) = \frac{1}{\sqrt{1-x^2}}$ дает пример, когда формула Чебышева *) имеет только действительные узлы при всяких n .

Были сделаны попытки строить формулы Чебышева для других весовых функций, но вычисления каждый раз показывали, что, начиная с некоторого значения n , среди узлов x_k будут существовать комплексные. Лишь сравнительно недавно были найдены весовые функции, для которых правило Чебышева (5.8.2) с действительными x_k может быть построено при всяких n , или для бесконечного числа значений n .

5.8.2. Случай постоянного веса $p(x) \equiv 1$

Отрезок интегрирования будем считать приведенным к $[-1, 1]$ и рассмотрим формулу

$$\int_{-1}^1 f(x)dx \approx C_n \sum_{k=1}^n f(x_k). \quad (5.8.6)$$

C_n и x_k нужно выбрать так, чтобы равенство было точным для степеней x от нулевой до n . Коэффициент C_n определится из условия, чтобы формула давала точный результат для $f \equiv 1$:

$$\int_{-1}^1 1 \cdot dx = 2 = C_n n, \quad C_n = \frac{2}{n}.$$

Ввиду $\int_{-1}^1 x^k dx = \frac{1}{k+1} [1 - (-1)^{k+1}]$, уравнения (5.8.3) для определения x_k здесь будут

*) Правило (5.6.8) является точным, когда f есть произвольный многочлен степени $2n-1$, а не только степени n , как это требуется для формулы Чебышева.

Это есть правило прямоугольников с высотой, равной ординате в средней точке.

$$\text{При } n=2 \quad \omega(x) = x^2 - \frac{1}{3}, \quad x_1 = -\frac{1}{\sqrt{3}}, \quad x_2 = \frac{1}{\sqrt{3}}, \quad C_2 = 1,$$

$$\int_{-1}^1 f(x) dx \approx f\left(-\frac{1}{\sqrt{3}}\right) + f\left(\frac{1}{\sqrt{3}}\right).$$

Правило выполняется точно для многочленов третьей степени и совпадает с формулой Ньютона для двух узлов.

$$\text{При } n=3 \quad \omega(x) = x^3 - \frac{1}{2}x, \quad x_1 = -\frac{1}{\sqrt{2}}, \quad x_2 = 0, \quad x_3 = \frac{1}{\sqrt{2}}, \quad C_3 = \frac{2}{3},$$

$$\int_{-1}^1 f(x) dx \approx \frac{2}{3} \left[f\left(-\frac{1}{\sqrt{2}}\right) + f(0) + f\left(\frac{1}{\sqrt{2}}\right) \right].$$

Приведем еще таблицу узлов формулы Чебышева для $n=1(1)7,9$.

$n=1$	$n=6$
$x_1=0$	$x_6=-x_1=0,86624\ 68181$
$n=2$	$x_5=-x_2=0,42251\ 86538$
$x_2=-x_1=0,57735\ 02691$	$x_4=-x_3=0,26663\ 54015$
$n=3$	$n=7$
$x_3=-x_1=0,70710\ 67812$	$x_7=-x_1=0,88386\ 17008$
$x_2=0$	$x_6=-x_2=0,52965\ 67753$
$n=4$	$x_5=-x_3=0,32391\ 18105$
$x_4=-x_1=0,79465\ 44723$	$x_4=0$
$x_3=-x_2=0,18759\ 24741$	$n=9$
$n=5$	$x_9=-x_1=0,91158\ 93007$
$x_5=-x_1=0,83249\ 74870$	$x_8=-x_2=0,60101\ 86554$
$x_4=-x_2=0,37454\ 14096$	$x_7=-x_3=0,52876\ 17831$
$x_3=0$	$x_6=-x_4=0,16790\ 61842$
	$x_5=0$

При $n=8$, как показали вычисления, среди x_k будут два комплексных. Расчеты были сделаны для нескольких $n>9$, но каждый раз оказывалось, что некоторые узлы x_k являются комплексными. В общем виде вопрос о возможности или невозможности построения правила Чебышева с действительными узлами для $n>9$ был решен в тридцатых годах текущего столетия С. Н. Бернштейном, показавшим, что при всяких $n>9$ среди узлов Чебышева будут комплексные.

Мы приведем доказательство теоремы Бернштейна, сохраняя порядок его рассуждений, но внеся в них некоторые упрощения. Докажем сначала несколько простых лемм.

Лемма 1. Допустим, что правило интегрирования

$$\int_{-1}^1 f(x) dx \approx \frac{2}{n} \sum_{k=1}^n f(x_k) \quad (5.8.9)$$

имеет действительные узлы $x_1 < x_2 < \dots < x_n$ и является точным для всякого многочлена степени $2m-1$, где $m < n$. Обозначим ξ_m наибольший корень многочлена Лежандра $P_m(x)$ степени m . Тогда $x_n > \xi_m$.

Доказательство. Положим $f = \frac{P_m^2(x)}{x - \xi_m} \cdot \frac{P_m(x)}{x - \xi_m}$ — это многочлен степени $m-1$, и он поэтому ортогонален к $P_m(x)$:

$$\int_{-1}^1 f dx = \int_{-1}^1 P_m(x) \frac{P_m(x)}{x - \xi_m} dx = 0.$$

$f(x)$ — это многочлен степени $2m-1$, для него равенство (5.8.9) должно быть точным и, стало быть,

$$\sum_{k=1}^n f(x_k) = 0.$$

Нулями $f(x)$ являются корни многочлена $P_m(x)$. Их m штук и, так как $m < n$, не все слагаемые $f(x_k)$ ($k=1, \dots, n$) равны нулю. Среди них должны быть как положительные, так и отрицательные. Но $f(x)$ принимает положительные значения при $x > \xi_m$ и отрицательные при $x < \xi_m$. Значит, для наибольшего узла непременно должно быть $x_n > \xi_m$.

Рассмотрим гауссово правило интегрирования с m узлами

$$\int_{-1}^1 f(x) dx \approx \sum_{i=1}^m A_i f(\xi_i), \quad (5.8.10)$$

$$P_m(\xi_i) = 0, \quad A_i = \frac{2}{(1 - \xi_i^2) [P_m'(\xi_i)]^2} \quad (i=1, 2, \dots, m).$$

Лемма 2. Если равенство (5.8.9) верно для всяких многочленов степени $2m-1$ ($m < n$), то

$$A_m > \frac{2}{n}. \quad (5.8.11)$$

Доказательство. Положим

$$f(x) = \left[\frac{P_m(x)}{(x-\xi_m)P'(\xi_m)} \right]^2.$$

f есть многочлен степени $2m-2$, обращающийся в единицу в узле ξ_m и в нуль в прочих узлах ξ_i ($i < m$). Квадратурная сумма Гаусса для f приводится к одному слагаемому $A_m f(\xi_m) = A_m \cdot 1$.

Правила (5.8.10) и (5.8.9) для f должны дать точное значение $\int_{-1}^1 f dx$, и поэтому

$$\frac{2}{n} \sum_{k=1}^n f(x_k) = A_m.$$

Так как $f(x) \geq 0$ при всяких x , отсюда следует

$$\frac{2}{n} f(x_n) \leq A_m. \quad (5.8.12)$$

Но из равенства

$$f(x) = [P_m'(\xi_m)]^{-2} (x-\xi_1)^2 \dots (x-\xi_{m-1})^2$$

видно, что $f(x)$ монотонно возрастает при $x \geq \xi_m$, и, ввиду $x_n > \xi_m$, будет

$$f(x_n) > f(\xi_m) = 1.$$

Отсюда и из (5.8.12) следует утверждение леммы.

Чтобы воспользоваться неравенством (5.8.11), нужно оценить

$$A_m = \frac{2}{(1-\xi_m^2) [P_m'(\xi_m)]^2},$$

для чего нам потребуется установить вспомогательные неравенства для ξ_m и $P_m'(\xi_m)$.

Лемма 3. При любом значении m для наибольшего корня ξ_m многочлена Лежандра $P_m(x)$ верно неравенство

$$1-\xi_m < \frac{3}{m(m+1)}. \quad (5.8.13)$$

Доказательство. Воспользуемся дифференциальным уравнением для $P_m(x)$:

$$\frac{d}{dx} [(1-x^2) P_m'(x)] + m(m+1) P_m(x) = 0.$$

После интегрирования его от ξ_m до 1 найдем

$$(1-\xi_m^2)P_m'(\xi_m) = m(m+1) \int_{\xi_m}^1 P_m(x) dx.$$

Если разложить $P_m(x)$ в ряд Тейлора по степеням $x-\xi_m$

$$P_m(x) = \sum_{i=1}^m \frac{(x-\xi_m)^i}{i!} P_m^{(i)}(\xi_m)$$

и выполнить почленное интегрирование, получим

$$(1-\xi_m^2)P_m'(\xi_m) = m(m+1) \sum_{i=1}^m \frac{(1-\xi_m)^{i+1}}{(i+1)!} P_m^{(i)}(\xi_m).$$

Корни многочлена $P_m(x)$ действительные, простые и лежат внутри отрезка $[-1, 1]$. По теореме Ролля, между каждыми двумя корнями ξ_i, ξ_{i+1} многочлена P_m лежит корень многочлена $P_m'(x)$. Стало быть, все корни P_m' являются простыми и лежат левее ξ_m . Аналогично, все корни P_m'' являются простыми и лежат левее ξ_m и т. д. Поэтому $P_m^{(i)}(\xi_m) > 0$ ($i=1, 2, \dots, m$) и все слагаемые в сумме последнего равенства положительны.

Сохраним в сумме два первых члена и отбросим остальные, кроме того, сократим обе части на положительный множитель $1-\xi_m$:

$$(1+\xi_m)P_m'(\xi_m) \geq m(m+1) \left[\frac{1}{2} (1-\xi_m)P_m'(\xi_m) + \frac{1}{6} (1-\xi_m)^2 P_m''(\xi_m) \right].$$

$P_m''(\xi_m)$ найдем из уравнения

$$(1-x^2)P_m'' - 2xP_m' + m(m+1)P_m = 0,$$

положив в нем $x = \xi_m$:

$$P_m''(\xi_m) = \frac{2\xi_m}{1-\xi_m^2} P_m'(\xi_m).$$

Подставим это значение в неравенство и сократим обе части на $P_m'(\xi_m)$:

$$1+\xi_m \geq m(m+1) \left[\frac{1}{2} (1-\xi_m) + \frac{1}{3} \frac{\xi_m(1-\xi_m)}{1+\xi_m} \right].$$

Усилим неравенство, заменив $1+\xi_m$ в знаменателе последнего члена правой части на большую величину 2:

$$1+\xi_m > m(m+1) \left[\frac{1}{2} (1-\xi_m) + \frac{1}{6} \xi_m(1+\xi_m) \right].$$

Отсюда, обозначив для сокращения $m(m+1) = \lambda$, получим:

$$\lambda \xi_m^2 + 2(3+\lambda)\xi_m + 6 - 3\lambda > 0. \quad (5.8.14)$$

Чтобы решить это неравенство, достаточно рассмотреть квадратное уравнение

$$\lambda z^2 + 2(3+\lambda)z + 6 - 3\lambda = 0.$$

$$z = \frac{\pm \sqrt{4\lambda^2 + 9 - 3 - \lambda}}{\lambda}.$$

Для нас имеет интерес лишь положительный корень z . Если ξ_m удовлетворяет неравенству (5.8.14), то ξ_m больше z :

$$\xi_m > \frac{\sqrt{4\lambda^2 + 9 - 3 - \lambda}}{\lambda} > \frac{2\lambda - 3 - \lambda}{\lambda} = 1 - \frac{3}{\lambda} = 1 - \frac{3}{m(m+1)}.$$

Лемма 3 доказана.

Лемма 4. Для значения производной многочлена Лежандра в наибольшем корне $x = \xi_m$ выполняется неравенство

$$P_m'(\xi_m) > \frac{2}{3(1-\xi_m)} \left[1 - \frac{\Gamma(m+4)}{288 \Gamma(m-2)} (1-\xi_m)^3 \right]. \quad (5.8.15)$$

Доказательство. Построим соотношение между ξ_m и $P_m'(\xi_m)$, которое для нас будет исходным. Запишем формулу Тейлора для $P_m(x)$ с остатком в виде интеграла:

$$P_m(x) = P_m'(\xi_m)(x-\xi_m) + \frac{1}{2} P_m''(\xi_m)(x-\xi_m)^2 + \frac{1}{2} \int_{\xi_m}^x P_m'''(t)(x-t)^2 dt.$$

Положим $x=1$ и вспомним, что $P_m(1)=1$:

$$1 = P_m'(\xi_m)(1-\xi_m) + \frac{1}{2} P_m''(\xi_m)(1-\xi_m)^2 + \frac{1}{2} \int_{\xi_m}^1 P_m'''(t)(1-t)^2 dt. \quad (5.8.16)$$

При доказательстве леммы 3 мы обращали внимание на то, что корни всех производных $P_m^{(i)}(x)$ ($i=1, 2, \dots, m-1$) лежат левее ξ_m . В частности, это верно для $P_m'''(x)$ и $P_m'''(x)$ будет монотонной возрастающей функцией на $[\xi_m, 1]$, достигающей своего наибольшего значения в точке $x=1$. Значение $P_m'''(1)$ можно при помощи несложных вычислений найти, если воспользоваться дифференциальным уравнением

$$(1-x^2)P_m''(x) - 2xP_m'(x) + m(m+1)P_m(x) = 0$$

и дважды его дифференцировать:

$$P_m'''(1) = \frac{\Gamma(m+4)}{48 \Gamma(m-2)}.$$

Если в (5.8.16) заменить $P_m''(\xi_m)$ ее значением, указанным выше, и $P_m'''(t)$ — большей величиной $P_m'''(1)$, получим неравенство

$$P_m'(\xi_m)(1-\xi_m) \left[1 + \frac{\xi_m}{1+\xi_m} \right] + \frac{\Gamma(m+4)}{48 \Gamma(m-2)} \cdot \frac{(1-\xi_m)^3}{6} \geq 1.$$

Так как $\frac{\xi_m}{1+\xi_m} < \frac{1}{2}$, отсюда сразу получаем (5.8.15).

На основании лемм 3 и 4 может быть построена нужная оценка для

$$A_m = \frac{2}{(1-\xi_m^2) [P_m'(\xi_m)]^2}.$$

$P_m'(\xi_m)$ заменим меньшей величиной из (5.8.15):

$$A_m < \frac{9(1-\xi_m)}{2(1+\xi_m)} \left[1 - \frac{\Gamma(m+4)}{288 \Gamma(m-2)} (1-\xi_m)^3 \right]^{-2}.$$

Для наших целей достаточна более грубая оценка. ξ_m возрастает при увеличении m . Будем рассматривать $m \geq 6$ и, так как $\xi_6 = 0,93246 \dots$, мы можем считать $1+\xi_m > 1,93$.

Далее, заменим везде $1-\xi_m$ большим числом $\frac{3}{m(m+1)}$.

Наконец, оценим величину, стоящую в квадратных скобках:

$$(m+3)(m-2) = m(m+1) - 6 < m(m+1),$$

$$(m+2)(m-1) = m(m+1) - 2 < m(m+1),$$

$$1 - \frac{\Gamma(m+4)}{288 \Gamma(m-2)} (1-\xi_m)^3 > 1 - \frac{m^3(m+1)^3}{288} \cdot \frac{3^3}{m^3(m+1)^3} = \frac{29}{32},$$

$$A_m < \frac{27 \cdot 32^2}{2 \cdot 1,93 \cdot 29^2} \cdot \frac{1}{m(m+1)} \approx \frac{8,517}{m(m+1)}. \quad (5.8.17)$$

Теорема 2. При $n \geq 10$ в правиле Чебышева (5.8.6) среди узлов x_k есть комплексные. Доказательство. Будем рассматривать такие значения n , при которых в правиле Чебышева (5.8.9) все узлы x_k действительные.

Предположим, что n есть число нечетное:

$$n = 2m - 1, \quad m = \frac{1}{2}(n+1), \quad m(m+1) = \frac{1}{4}(n+1)(n+3).$$

Для A_m должно выполняться неравенство

$$A_m < \frac{4 \cdot 8,517}{(n+1)(n+3)},$$

Ввиду леммы 2, должно также быть

$$\frac{4 \cdot 8,517}{(n+1)(n+3)} > \frac{2}{n}$$

или

$$n^2 - 13,034n + 3 < 0, \quad n < 13.$$

Следовательно, при $n \geq 13$ в правиле Чебышева не все x_k являются действительными. Для $n=11$ все x_k также не могут быть действительными, так как тогда $m=6$, $A_6 = 0,173 \dots$, $\frac{2}{11} = 0,1818$ и неравенство $\frac{2}{11} < A_6$ не выполняется.

Пусть n — четное. Правило Чебышева (5.8.6) будет точным для многочленов степени $n+1$. Положим

$$n+1 = 2m-1, \quad m = \frac{1}{2}(n+2).$$

Согласно (5.8.11) и (5.8.17), должно быть

$$\frac{4 \cdot 8,517}{(n+2)(n+4)} > \frac{2}{n} \quad \text{и} \quad n < 11.$$

Стало быть, для четных $n > 10$ среди узлов Чебышева есть комплексные. При $n=10$ это также верно, так как неравенство $A_6 = 0,173 \dots > \frac{2}{10} = 0,2$ не выполняется.

§ 5.9. УВЕЛИЧЕНИЕ ТОЧНОСТИ КВАДРАТУРНЫХ ПРАВИЛ. ФОРМУЛЫ ЭЙЛЕРОВА ВИДА

5.9.1. Введение

Выберем какое-либо определенное квадратурное правило и рассмотрим его остаток

$$R(f) = \int_a^b p(x)f(x)dx - \sum_{k=1}^n A_k f(x_k).$$

Численное значение его зависит от двух фактов: от свойств интегрируемой функции f и от свойств избранного правила.

Поэтому в задаче уменьшения погрешности вычисления интеграла можно наметить два направления.

1. Правила приближенных квадратур, о которых говорилось выше, были построены при помощи замены интегрируемой функции f алгебраи-

ческим многочленом и рассчитаны, следовательно, на интегрирование функций, имеющих тот или иной порядок гладкости. Так, например, погрешность правила Гаусса с n узлами зависит от того, насколько точно f может быть на всем отрезке $\langle a, b \rangle$ приближена многочленом степени $2n-1$. Аналогично, погрешность правила трапеций (5.4.3) зависит от того, как сильно на каждом частичном отрезке $[a+kh, a+(k+1)h]$ график f будет отличаться от прямолинейной хорды, соединяющей концы соответствующего участка графика.

Если функция f в недостаточной степени обладает теми свойствами, при которых можно ожидать хорошей точности результата, например будет разрывной или непрерывной, но имеющей разрывную первую производную, она не может быть хорошо приближена многочленом невысокой степени, и тогда трудно ожидать малой погрешности $R(f)$ при ее приближенном интегрировании. Иногда большое значение погрешности можно получить при интегрировании аналитических функций, если их особенности лежат вблизи отрезка $\langle a, b \rangle$. В этих случаях полезно перед интегрированием предварительно преобразовать функцию f так, чтобы устранить или ослабить те ее свойства, которые могут вызвать большие значения $R(f)$.

Нередко можно значительно увеличить точность вычисления, если выделить из f «особую часть» путем разложения f на два слагаемых $f = f_1 + f_2$ так, чтобы f_1 содержала «все» особенности f или «главную часть» особенностей. Кроме того, f_1 должна быть такой, чтобы интеграл $\int_a^b p f_1 dx$ вычислялся точно. Второе же слагаемое f_2 либо совсем не должно иметь особенностей, либо его особенности должны быть настолько слабее особенностей f , чтобы интеграл $\int_a^b p f_2 dx$ мог быть вычислен при помощи взятого правила приближенной квадратуры с достаточной точностью. Некоторые способы выделения и ослабления особенностей функций будут рассмотрены в следующем параграфе.

2. Рассмотрим ту часть погрешности, которая вызвана недостаточной точностью избранного правила. Чтобы увеличить точность правила, нужно к квадратурной сумме $\sum_{k=1}^n A_k f(x_k)$ добавить дополнительное слагаемое и выбрать его так, чтобы оно являлось главной частью остатка $R(f)$. От него мы должны потребовать вычислимости и достаточной простоты.

Один из методов выделения главной части из $R(f)$, приводящий к естественному обобщению формулы Эйлера — Маклорена, будет изложен в следующем пункте.

Пусть новый член квадратурного правила найден. Если добавление его к квадратурной сумме исправит ранее найденный результат до нужной точности, на этом улучшение правила можно закончить. Если же

нужная точность не будет достигнута, то следует найти остаток улучшенного правила и выделить из него в свою очередь главную часть и т. д.

В различных задачах приходится выделять разное число главных частей и при построении теории уточнения правил мы должны предусмотреть разложение $R(f)$ в ряд, состоящий из главных частей возрастающих порядков, и найти остаток, получающийся после выделения из $R(f)$ любого конечного числа членов такого разложения.

5.9.2. Правила эйлерова вида

Рассмотрим квадратурное правило

$$\int_a^b p(x)f(x)dx = \sum_{k=1}^n A_k f(x_k) + R(f) \quad (5.9.1)$$

и предположим, что оно имеет степень точности $m-1$. Отрезок интегрирования $[a, b]$ считается конечным. Остаток $R(f)$, если воспользоваться тейлоровым разложением функции

$$f(x) = \sum_{i=1}^{m-1} \frac{(x-a)^i}{i!} f^{(i)}(a) + \int_a^b f^{(m)}(t) E(x-t) \frac{(x-t)^{m-1}}{(m-1)!} dt$$

и внести его в равенство

$$R(f) = \int_a^b p f dx - \sum_{k=1}^n A_k f(x_k),$$

как это мы делали в конце § 5.2, можно представить в форме, достаточно удобной для наших целей:

$$R(f) = \int_a^b f^{(m)}(t) K(t) dt, \quad (5.9.2)$$

$$K(t) = \int_t^b p(x) \frac{(x-t)^{m-1}}{(m-1)!} dx - \sum_{k=1}^n A_k E(x_k - t) \frac{(x_k - t)^{m-1}}{(m-1)!}.$$

Способы выделения из $R(f)$ главной части тесно связаны со свойствами ядра $K(t)$. Если значения ядра в какой-то мере равномерно распределены вдоль всего отрезка $[a, b]$ и ядро не сильно изменяется, то на образование $R(f)$ наибольшее влияние будет оказывать среднее зна-

чение ядра. Чтобы пояснить эту мысль примером, возьмем общую формулу трапеции (5.4.3). Она получена путем деления отрезка $[a, b]$ на некоторое число одинаковых частей «малой длины» h и применения к каждой части элементарного правила трапеций (5.4.1). Формула верна для линейной функции и остаток ее выражается через вторую производную f'' . Совершенно ясно, что если для остатка формулы (5.4.3) написать интегральное представление типа (5.9.2), оно будет иметь форму

$$R(f) = \int_a^b f''(t) K(t) dt$$

и ядро $K(t)$ его будет на $[a, b]$ периодической функцией с «малым периодом» h . Оно будет обладать указанным выше свойством.

Аналогичное можно сказать об общих правилах парабол (5.4.6) и «трех восьмых» (5.4.9), так же как о всяком другом правиле приближенного интегрирования, основанном на делении $[a, b]$ на «малые части», применении к каждой части какого-либо одного и того же правила и последующем сложении результатов.

Полученные ниже результаты верны для всякого квадратурного правила (5.9.1). Наглядные соображения, которые были приведены выше, позволяют лишь предвидеть, для каких правил полученные формулы могут дать хорошие результаты при применении их для улучшения точности вычисления интеграла.

Для выделения из $R(f)$ (5.9.2) главной части достаточно положить

$$\begin{aligned} C_0 &= (b-a)^{-1} \int_a^b K(t) dt, \quad K(t) = C_0 + [K(t) - C_0], \\ R(f) &= C_0 \int_a^b f^{(m)}(t) dt - \int_a^b f^{(m)}(t) [C_0 - K(t)] dt = C_0 [f^{(m-1)}(b) - f^{(m-1)}(a)] + \\ &+ \int_a^b f^{(m+1)}(t) L_1(t) dt, \quad L_1(t) = \int_a^t [C_0 - K(x)] dx. \end{aligned}$$

При получении последней части системы равенств в интеграле

$$\int_a^b f^{(m)} [C_0 - K] dt$$

было выполнено интегрирование по частям.

Из интеграла

$$\int_a^b f^{(m+1)}(t) L_1(t) dt$$

в свою очередь может быть выделена главная часть и т. д.

После s -кратного выделения из остатка главных частей получим для улучшения точности формулы (5.9.1) следующее правило эйлера вида:

$$\begin{aligned} \int_a^b p(x) f(x) dx = & \sum_{k=1}^n A_k f(x_k) + C_0 [f^{(m-1)}(b) - f^{(m-1)}(a)] + \dots + \\ & + C_{s-1} [f^{(m+s-2)}(b) - f^{(m+s-2)}(a)] + R_s(f), \end{aligned} \quad (5.9.3)$$

$$C_i = (b-a)^{-1} \int_a^b L_i(t) dt, \quad L_{i+1}(t) = \int_a^t [C_i - L_i(x)] dx, \quad L_0(t) = K(t),$$

$$R_s(f) = \int_a^b f^{(m+s)}(t) L_s(t) dt. \quad (5.9.4)$$

Исходное правило (5.9.1) было, по предположению, точным для многочленов степени $m-1$. Если к квадратурной сумме $\sum_{k=1}^n A_k f(x_k)$ прибавить слагаемое $C_0 [f^{(m-1)}(b) - f^{(m-1)}(a)]$, получится правило, верное для многочленов степени m , добавление еще второго слагаемого $C_1 [f^{(m)}(b) - f^{(m)}(a)]$ сделает правило верным для многочленов степени $m+1$ и т. д.

Формулы (5.9.4) позволяют находить последовательно C_i и $L_i(t)$. Можно построить их выражения непосредственно через ядро $K(t)$. Для этого в интегральном выражении (5.9.2) для $K(t)$ заменим $f^{(m)}(t)$ ее разложением по многочленам Бернулли *) [добавление II, (II.21)]

$$\begin{aligned} f^{(m)}(t) = & (b-a)^{-1} \int_a^b f^{(m)}(x) dx + \\ & + \sum_{i=1}^{s-1} \frac{(b-a)^{i-1}}{i!} B_i \left(\frac{t-a}{b-a} \right) [f^{(m+i-1)}(b) - f^{(m+i-1)}(a)] - \end{aligned}$$

) Здесь $B_i(x)$ есть многочлен Бернулли степени i и $B_i^(x)$ есть 1-периодическая функция, совпадающая с $B_i(x)$ на промежутке $0 \leq x < 1$.

$$- \frac{(b-a)^{s-1}}{s!} \int_a^b f^{(m+s)}(x) \left[B_s^* \left(\frac{t-x}{b-a} \right) - B_s^* \left(\frac{t-a}{b-a} \right) \right] dx.$$

Это приведет к равенству

$$\begin{aligned} \int_a^b p(x) f(x) dx = & \sum_{k=1}^n A_k f(x_k) + (b-a)^{-1} \int_a^b K(t) dt [f^{(m-1)}(b) - f^{(m-1)}(a)] + \\ & + \sum_{i=1}^{s-1} \frac{(b-a)^{i-1}}{i!} \int_a^b K(t) B_i \left(\frac{t-a}{b-a} \right) dt [f^{(m+i-1)}(b) - f^{(m+i-1)}(a)] - \\ & - \frac{(b-a)^{s-1}}{s!} \int_a^b K(t) \int_a^b f^{(m+s)}(x) \left[B_s^* \left(\frac{t-x}{b-a} \right) - B_s^* \left(\frac{t-a}{b-a} \right) \right] dx dt. \end{aligned}$$

Оно должно совпадать с (5.9.3) для всякой функции f , имеющей на $[a, b]$ непрерывную производную порядка $m+s$, что может быть только в том случае, когда будут одинаковыми коэффициенты при $f^{(m+i-1)}(b) - f^{(m+i-1)}(a)$ ($i=1, 2, \dots, s-1$) и множители при $f^{(m)}(t)$ в интегральных выражениях для остаточных членов:

$$\left. \begin{aligned} C_i &= \frac{(b-a)^{i-1}}{i!} \int_a^b K(t) B_i \left(\frac{t-a}{b-a} \right) dt, \\ L_s(t) &= - \frac{(b-a)^{s-1}}{s!} \int_a^b K(x) \left[B_s^* \left(\frac{x-t}{b-a} \right) - B_s^* \left(\frac{x-a}{b-a} \right) \right] dx. \end{aligned} \right\} \quad (5.9.5)$$

Сделаем добавление к полученному результату. Из сравнения выражения (5.9.5) для C_i с представлением (5.9.2) остатка видно, что C_i есть остаток квадратуры некоторой функции, производная которой порядка m равна

$$\frac{(b-a)^{i-1}}{i!} B_i \left(\frac{t-a}{b-a} \right).$$

Правило дифференцирования многочленов Бернулли [добавление II, (II. 11)] говорит, что за такую функцию можно принять

$$\frac{(b-a)^{m+i-1}}{(m+i)!} B_{m+i} \left(\frac{t-a}{b-a} \right)$$

и, стало быть,

$$C_i = \frac{(b-a)^{m+i-1}}{(m+i)!} R \left[B_{m+i} \left(\frac{t-a}{b-a} \right) \right] =$$

$$= \frac{(b-a)^{m+i-1}}{(m+i)!} \left\{ \int_a^b p(t) B_{m+i} \left(\frac{t-a}{b-a} \right) dt - \sum_{k=1}^n A_k B_{m+i} \left(\frac{x_k-a}{b-a} \right) \right\}. \quad (5.9.6)$$

Сходные соображения применимы к $L_s(t)$ и

$$L_s(t) = - \frac{(b-a)^{m+s-1}}{(m+s)!} R_x \left[B_{m+s}^* \left(\frac{x-t}{b-a} \right) - B_{m+s}^* \left(\frac{x-a}{b-a} \right) \right], \quad (5.9.7)$$

где знак x , стоящий около R , показывает, что вычисляется остаток квадратуры по переменной x , тогда как величина t является параметром.

5.9.3. Формула Эйлера — Маклорена

Рассмотрим простейшее правило трапеций и построим для него уточняющее равенство (5.9.3).

$$\int_a^b f(x) dx = \frac{b-a}{2} [f(a) + f(b)] + R(f). \quad (5.9.8)$$

Алгебраическая степень точности правила равна единице, и нужно считать $m=2$. Для вычисления C_i воспользуемся (5.9.6). Вспомним, что многочлены Бернулли $B_n(z)$ ($n=2, 3, \dots$) принимают в точках $z=0$ и $z=1$ одинаковые значения и, следовательно,

$$\int_a^b B_{i+2} \left(\frac{t-a}{b-a} \right) dt = \frac{b-a}{i+3} [B_{i+3}(1) + B_{i+3}(0)] = 0,$$

$$C_i = - \frac{(b-a)^{i+2}}{(i+2)!} \cdot \frac{1}{2} [B_{i+2}(0) + B_{i+2}(1)] =$$

$$= - \frac{(b-a)^{i+2}}{(i+2)!} \cdot \frac{1 + (-1)^{i+2}}{2!} B_{i+2}.$$

Отсюда видно, что при нечетных значениях i все C_i равны нулю: $C_1 = C_3 = C_5 = \dots = 0$. Для четных значений $i=2k$ будет:

$$C_{2k} = - \frac{(b-a)^{2k+2}}{(2k+2)!} B_{2k+2}.$$

Остаток $R(f)$ формулы (5.9.8) найдем, воспользовавшись (5.9.4) и (5.9.7):

$$L_s(t) = -\frac{(b-a)^{s+1}}{(s+2)!} \left\{ \int_a^b \left[B_{s+2}^* \left(\frac{x-t}{b-a} \right) - B_{s+2}^* \left(\frac{x-a}{b-a} \right) \right] dx - \right. \\ \left. - \frac{b-a}{2} \left[\left(B_{s+2}^* \left(\frac{a-t}{b-a} \right) - B_{s+2}^*(0) \right) + \left(B_{s+2}^* \left(\frac{b-t}{b-a} \right) - B_{s+2}^*(1) \right) \right] \right\}.$$

Так как

$$B_{s+2}^* \left(\frac{x-t}{b-a} \right)$$

является $(b-a)$ -периодической функцией x , то интегралы

$$\int_a^b B_{s+2}^* \left(\frac{x-t}{b-a} \right) dx \quad \text{и} \quad \int_a^b B_{s+2}^* \left(\frac{x-a}{b-a} \right) dx$$

имеют одно и то же значение и интегральный член в $L_s(t)$ исчезает. Далее,

$$B_{s+2}^* \left(\frac{a-t}{b-a} \right) = B_{s+2}^* \left(\frac{b-t}{b-a} \right), \quad B_{s+2}^*(0) = B_{s+2}^*(1) = B_{s+2}$$

и, следовательно,

$$L_s(t) = \frac{(b-a)^{s+2}}{(s+2)!} \left[B_{s+2}^* \left(\frac{b-t}{b-a} \right) - B_{s+2} \right] = \frac{(b-a)^{s+2}}{(s+2)!} y_{s+2}^* \left(\frac{b-t}{b-a} \right).$$

Для простейшего правила трапеций (5.9.8) можно теперь образовать разложение (5.9.3). Все члены, содержащие C_i нечетных индексов, будут отсутствовать. Считая s числом четным и полагая $s+2=2v$, можно записать разложение (5.9.3) в форме

$$\int_a^b f(t) dt = \frac{b-a}{2} [f(a) + f(b)] - \\ - \sum_{k=1}^{v-1} \frac{(b-a)^{2k}}{(2k)!} B_{2k} [f^{(2k-1)}(b) - f^{(2k-1)}(a)] + \rho_{2v}(f), \quad (5.9.9) \\ \rho_{2v}(f) = \frac{(b-a)^{2v}}{(2v)!} \int_a^b f^{(2v)}(t) y_{2v}^* \left(\frac{b-t}{b-a} \right) dt.$$

Ниже нам удобнее будет пользоваться другой формой остаточного члена $\rho_{2v}(f)$. Положим $t = a + (b-a)u$ ($0 \leq u \leq 1$). Так как

$$y_{2v}^* \left(\frac{b-t}{b-a} \right) = y_{2v}^*(1-u) = B_{2v}(1-u) - B_{2v} = B_{2v}(u) - B_{2v} = y_{2v}(u),$$

для остатка найдем

$$\begin{aligned} \rho_{2v}(f) &= \frac{(b-a)^{2v+1}}{(2v)!} \int_0^1 f^{(2v)}[a + (b-a)u] [B_{2v}(u) - B_{2v}] du = \\ &= -\frac{(b-a)^{2v+1}}{(2v)!} \int_0^1 f^{(2v)}[a + (b-a)u] y_{2v}(u) du. \end{aligned} \quad (5.9.10)$$

Для получения правила увеличения точности общей формулы трапеций (5.4.3) отрезок $[a, b]$ разделим на n одинаковых частей точками $a + ph$ ($p=0, 1, \dots, n$), $h = \frac{1}{n}(b-a)$ и применим равенство (5.9.9) к частичному отрезку $[a + ph, a + (p+1)h]$.

$$\begin{aligned} \int_{a+ph}^{a+(p+1)h} f(x) dx &= \frac{h}{2} \{f[a+ph] + f[a+(p+1)h]\} - \\ &- \sum_{k=1}^{v-1} \frac{h^{2k}}{(2k)!} B_{2k} \{f^{(2k-1)}[a+(p+1)h] - f^{(2k-1)}[a+ph]\} + \rho_{2v}^{(p)}(f), \\ \rho_{2v}^{(p)}(f) &= \frac{h^{2v+1}}{(2v)!} \int_0^1 f^{(2v)}[a+h(p+u)] y_{2v}(u) du. \end{aligned}$$

Если суммировать такие равенства для всех отрезков ($p=0, 1, \dots, n-1$), слагаемые в суммах $\sum_{k=1}^{v-1}$, отвечающие точкам деления, которые лежат внутри $[a, b]$, сократятся и останутся лишь слагаемые, отвечающие концам a и b отрезка интегрирования, и мы получим широко известную формулу Эйлера — Маклорена

$$\int_a^b f(x) dx = T_n - \sum_{k=1}^{v-1} \frac{h^{2k}}{(2k)!} B_{2k} [f^{(2k-1)}(b) - f^{(2k-1)}(a)] + \rho_{2v}(f) =$$

$$\begin{aligned}
&= T_n - \frac{h^2}{12} [f'(b) - f'(a)] + \frac{h^4}{720} [f'''(b) - f'''(a)] - \\
&- \frac{h^6}{30\,240} [f^{(5)}(b) - f^{(5)}(a)] + \frac{h^8}{1\,209\,600} [f^{(7)}(b) - f^{(7)}(a)] - \\
&- \frac{h^{10}}{47\,900\,160} [f^{(9)}(b) - f^{(9)}(a)] + \dots + \rho_{2v}(f), \quad (5.9.11)
\end{aligned}$$

где

$$T_n = h \left[\frac{1}{2} f(a) + f(a+h) + f(a+2h) + \dots + \frac{1}{2} f(b) \right],$$

$$\rho_{2v}(f) = \frac{h^{2v+1}}{(2v)!} \int_0^1 y_{2v}(u) \sum_{p=0}^{n-1} f^{(2v)}(a+ph+uh) du.$$

Если неограниченно увеличивать v , то сумма $\sum_{k=1}^{v-1}$ в пределе дает ряд

$$\sum_{k=1}^{\infty} \frac{h^{2k}}{(2k)!} B_{2k} [f^{(2k-1)}(b) - f^{(2k-1)}(a)].$$

В добавлении II показывается, что числа Бернулли B_{2k} с ростом k начинают быстро возрастать, при больших k верно приближенное равенство

$$B_{2k} \approx 2(-1)^{k-1} (2k)! (2\pi)^{-2k}$$

и ряд будет сходиться для узкого множества функций f . В широком классе случаев члены ряда, начиная с некоторого номера, неограниченно возрастают и формула Эйлера — Маклорена не даст возможности вычислить интеграл сколь угодно точно. Но часто оказывается, что остаток $\rho_{2v}(f)$ для первых значений v убывает с ростом v и имеет малую величину, и если это выполняется, то формула (5.9.11) может принести заметную пользу в уточнении правила трапеций.

При изучении остатка $\rho_{2v}(f)$ нередко можно воспользоваться приводимыми ниже двумя теоремами.

Теорема 1. Если $f^{(2v)}(x)$ непрерывна на отрезке $[a, b]$, то существует такая точка ξ ($a \leq \xi \leq b$), что для остатка $\rho_{2v}(f)$ верно равенство

$$\rho_{2v}(f) = - \frac{h^{2v+1}}{(2v)!} B_{2v} f^{(2v)}(\xi) \quad (5.9.12)$$

Доказательство. Для доказательства рассмотрим интеграл

$$I = \int_0^1 y_{2\nu}(u) \sum_{p=0}^{n-1} f^{(2\nu)}(a+ph+hu) du = \int_0^1 y_{2\nu}(u) \sigma(u) du.$$

В добавлении II показывается, что $y_{2\nu}(x)$ сохраняет знак на отрезке $0 < u < 1$ и здесь применима теорема о среднем взвешенном значении:

$$I = \sigma(v) \int_0^1 y_{2\nu}(u) du \quad (0 \leq v \leq 1).$$

Если M и m есть наибольшее и наименьшее значения $f^{(2\nu)}(u)$ ($0 \leq u \leq 1$), то $\sigma(u)$, очевидно, лежит в следующих границах: $nm \leq \sigma(u) \leq nM$. Поэтому для $\sigma(v)$ имеет место равенство $\sigma(v) = nP$ при $m \leq P \leq M$. Ввиду же непрерывности $f^{(2\nu)}(u)$, существует на $[a, b]$ такая точка ξ , что $\sigma(v) = nf^{(2\nu)}(\xi)$. Кроме того,

$$\int_0^1 y_{2\nu}(u) du = \int_0^1 [B_{2\nu}(u) - B_{2\nu}] du = -B_{2\nu}$$

и, значит, $I = -B_{2\nu}nf^{(2\nu)}(\xi)$. Отсюда и из выражения остатка $\rho_{2\nu}(f)$, указанного в (5.9.11), сразу следует утверждение теоремы.

Теорема 2. Если $f^{(2\nu+2)}(x)$ непрерывна на $[a, b]$ и при всяких x ($a \leq x \leq b$) будет либо $f^{(2\nu)}(x) \geq 0$ и $f^{(2\nu+2)}(x) \geq 0$, либо $f^{(2\nu)}(x) \leq 0$ и $f^{(2\nu+2)}(x) \leq 0$, то величины $\rho_{2\nu}(f)$ и $-\rho_{2\nu+2}(f)$ имеют такие же знаки, как

$-\frac{h^{2\nu}}{(2\nu)!} B_{2\nu} [f^{(2\nu-1)}(b) - f^{(2\nu-1)}(a)]$ и по абсолютному значению не больше этого члена формулы.

Доказательство. Для остатков $\rho_{2\nu}(f)$ и $\rho_{2\nu+2}(f)$ верна следующая связь:

$$\rho_{2\nu}(f) = -\frac{h^{2\nu}}{(2\nu)!} B_{2\nu} [f^{(2\nu-1)}(b) - f^{(2\nu-1)}(a)] + \rho_{2\nu+2}(f).$$

Ее можно переписать в виде

$$\rho_{2\nu}(f) + [-\rho_{2\nu+2}(f)] = \frac{h^{2\nu+1}}{(2\nu)!} \int_0^1 y_{2\nu}(u) \sum_{p=0}^{n-1} f^{(2\nu)}(a+ph+hu) du +$$

$$\begin{aligned}
& + \frac{h^{2v+3}}{(2v+2)!} \int_0^1 [-y_{2v+2}(u)] \sum_{p=0}^{n-1} f^{(2v+2)}(a+ph+hu) du = \\
& = - \frac{h^{2v}}{(2v)!} B_{2v} [f^{(2v-1)}(b) - f^{(2v-1)}(a)].
\end{aligned}$$

По предположению, $f^{(2v)}(x)$ и $f^{(2v+2)}(x)$ сохраняют одинаковые знаки всюду на $[a, b]$. В добавлении II доказывается, что на $[0, 1]$ $y_{2v}(u)$ и $-y_{2v+2}(u)$ также сохраняют один и тот же знак [добавление II, § 2]. Поэтому $\rho_{2v}(f)$ и $-\rho_{2v+2}(f)$ имеют один и тот же знак. Этот знак должен совпадать со знаком последней части, и каждая из этих величин не больше по модулю последнего члена равенства.

Формула Эйлера — Маклорена является единственным конкретным правилом эйлерова вида для увеличения точности механических квадратур, на котором мы остановились. Но такие правила, как отмечалось выше, могут быть построены для каждой квадратурной формулы.

Для ознакомления с некоторыми из них мы отсылаем к справочной литературе [4].

5.9.4. Разностные видоизменения формулы Эйлера — Маклорена

Формула Эйлера — Маклорена требует вычисления производных f', f'', \dots на концах отрезка интегрирования, что не всегда просто и даже не всегда возможно. Можно построить несколько разновидностей этой формулы, в которых уточняющие члены выражаются через значения функции и не содержат производных. Все они могут быть получены из (5.9.11) путем замены там производных на приближенные выражения их через значения f в точках $a+kh$ ($k=0, \pm 1, \pm 2, \dots$). Замена может быть сделана многими способами и можно построить не одно, а несколько таких правил. Мы остановимся только на одном из них, в котором используются узлы, не выходящие за границу отрезка $[a, b]$.

Начнем с вычисления производных в точке a . Интерполируем $f(x)$ по ее значениям в точках $a, a+h, a+2h, \dots$. Это может быть сделано при помощи формулы Ньютона (4.4.1) для интерполирования в начале таблицы

$$f(x) = f(a+th) = f_0 + \frac{t}{1!} \Delta f_0 + \frac{t(t-1)}{2!} \Delta^2 f_0 + \dots + r(x), \quad f_k = f(a+kh).$$

Вычислив производные и полагая $x=a, t=0$, получим: *)

*) Равенства были получены в конце § 4.5 и воспроизведены здесь для облегчения чтения.

$$\left. \begin{aligned} hf'(a) &= \Delta f_0 - \frac{1}{2} \Delta^2 f_0 + \frac{1}{3} \Delta^3 f_0 - \frac{1}{4} \Delta^4 f_0 + \frac{1}{5} \Delta^5 f_0 + \dots + r'(a), \\ h^2 f''(a) &= \Delta^2 f_0 - \Delta^3 f_0 + \frac{11}{12} \Delta^4 f_0 - \frac{5}{6} \Delta^5 f_0 + \dots + r''(a), \\ h^3 f'''(a) &= \Delta^3 f_0 - \frac{3}{2} \Delta^4 f_0 + \frac{7}{4} \Delta^5 f_0 - \dots + r'''(a), \\ h^4 f^{(4)}(a) &= \Delta^4 f_0 - 2\Delta^5 f_0 + \dots, \\ h^5 f^{(5)}(a) &= \Delta^5 f_0 - \dots \end{aligned} \right\} \quad (5.9.13)$$

Для нахождения производных на правом конце $x=b$ выполним интерполирование f по значениям в точках $b=a+nh$, $a+(n-1)h$, $a+(n-2)h$, ..., воспользовавшись формулой Ньютона (4.4.3) для интерполирования в конце таблицы

$$f(x) = f(b+th) = f_n + \frac{t}{1!} \Delta f_{n-1} + \frac{t(t+1)}{2!} \Delta^2 f_{n-2} + \dots + \rho(x).$$

Отсюда при $x=b=a+nh$, $t=0$ получим:

$$\left. \begin{aligned} hf'(b) &= \Delta f_{n-1} + \frac{1}{2} \Delta^2 f_{n-2} + \frac{1}{3} \Delta^3 f_{n-3} + \frac{1}{4} \Delta^4 f_{n-4} + \frac{1}{5} \Delta^5 f_{n-5} + \dots + \rho'(b), \\ h^2 f''(b) &= \Delta^2 f_{n-2} + \Delta^3 f_{n-3} + \frac{11}{12} \Delta^4 f_{n-4} + \frac{5}{6} \Delta^5 f_{n-5} + \dots, \\ h^3 f'''(b) &= \Delta^3 f_{n-3} + \frac{3}{2} \Delta^4 f_{n-4} + \frac{7}{4} \Delta^5 f_{n-5} + \dots, \\ h^4 f^{(4)}(b) &= \Delta^4 f_{n-4} + 2\Delta^5 f_{n-5} + \dots, \\ h^5 f^{(5)}(b) &= \Delta^5 f_{n-5} + \dots \end{aligned} \right\} \quad (5.9.14)$$

Отбросим теперь остаточные члены в равенствах (5.9.13) и (5.9.14) и полученные приближенные значения производных внесем в формулу Эйлера — Маклоренга. После некоторых преобразований получится формула Грегори

$$\begin{aligned} \int_a^{a+nh} f(x) dx &= T_n - \frac{h}{12} (\Delta f_{n-1} - \Delta f_0) - \frac{h}{24} (\Delta^2 f_{n-2} + \Delta^2 f_0) - \\ &- \frac{19h}{720} (\Delta^3 f_{n-3} - \Delta^3 f_0) - \frac{3h}{160} (\Delta^4 f_{n-4} + \Delta^4 f_0) - \frac{863h}{60 \cdot 480} (\Delta^5 f_{n-5} - \Delta^5 f_0) - \end{aligned}$$

$$-\frac{275h}{24 \cdot 192} (\Delta^6 f_{n-6} + \Delta^6 f_0) - \dots - C_k h [\Delta^k f_{n-k} + (-1)^k \Delta^k f_0] + R_1(f), \quad (5.9.15)$$

$$C_k = \frac{(-1)^k}{(k+1)!} \int_0^1 x(x-1) \dots (x-k) dx.$$

§ 5.10. УВЕЛИЧЕНИЕ ТОЧНОСТИ КВАДРАТУРНЫХ ПРАВИЛ. ОСЛАБЛЕНИЕ ОСОБЕННОСТЕЙ ИНТЕГРИРУЕМОЙ ФУНКЦИИ

В настоящем параграфе мы рассмотрим некоторые вопросы, связанные с подготовкой функции к интегрированию, и укажем несколько практически полезных для этого правил. Чтобы выяснить идеи, положенные в основу предварительных преобразований, достаточно рассмотреть интеграл простейшего вида с постоянной весовой функцией и конечным отрезком интегрирования.

Каждое правило приближенной квадратуры, рассмотренное нами выше,

$$\int_a^b f(x) dx = \sum_{k=1}^n A_k f(x_k) + R(f), \quad (5.10.1)$$

было основано на замене функции f алгебраическим многочленом на всем отрезке $\langle a, b \rangle$ или на его частях, и следует ожидать, что правило может дать хорошую точность, если f обладает «достаточно высоким порядком» гладкости. Поэтому одной из первых целей, которую ставят при предварительной подготовке функции к интегрированию, является повышение ее гладкости до границ, оптимальных для избранного правила.

За меру гладкости функции принимают порядок ее непрерывной дифференцируемости и улучшение гладкости функции означает в первую очередь повышение этого порядка. Достигается это путем выделения из f ее «особой части». С некоторыми правилами такого выделения мы ознакомимся ниже. Отметим, что повышение порядка дифференцируемости, хотя и является весьма полезным для улучшения точности, имеет ограниченное значение. Границы его целесообразности для каждого правила будут свои и зависят от степени точности правила. Пояснить это обстоятельство можно наиболее просто на примере. Рассмотрим общее правило парабол (5.4.6). Оно дает точный результат, если f есть многочлен третьей степени. Остаток его может быть выражен через четвертую производную функции f . Для $R(f)$ было получено представление (5.4.7)

$$R = -\frac{(b-a)^5}{180n^4} f^{IV}(\xi) \quad (a \leq \xi \leq b).$$

Из него, в частности, видно, что $R \neq 0$, если $f^{IV}(x)$ отлична от нуля всюду на $[a, b]$. Нам удобнее воспользоваться другим представлением остатка, не содержащим неизвестных величин. Так как степень точности правила парабол равна 3, для ее остатка будет справедлива формула (5.9.2) при $m=4$:

$$R(f) = \int_a^b f^{IV}(x) K(x) dx. \quad (5.10.2)$$

Явное выражение для ядра $K(x)$ интеграла нам сейчас не потребуется, и мы не станем его приводить. Но полезно отметить, что ядро $K(x)$ непрерывно и не изменяет знака на $[a, b]$, так как если бы $K(x)$ меняло знак, то существовала бы такая функция f , имеющая непрерывную и сохраняющую знак производную f^{IV} , для которой $R(f) = 0$, а это, как отмечалось выше, невозможно.

Возвратимся к задаче об увеличении порядка дифференцируемости f . При применении правила парабол естественно стремиться к тому, чтобы интегрируемая функция f была трижды или, лучше, четырежды непрерывно дифференцируемой. Пусть мы достигли четырехкратной дифференцируемости f . Нас интересует вопрос, можно ли ожидать значительного увеличения точности, если стремиться к дальнейшему повышению порядка дифференцируемости f . Ответ на этот вопрос может быть получен при первом взгляде на (5.10.2). Для всякой функции сколь угодно высокого порядка дифференцируемости, лишь бы он не был меньше четырех, остаток всегда представим в форме (5.10.2) и при одном только повышении порядка дифференцируемости, без целесообразного изменения свойств четвертой производной, уменьшение остатка $R(f)$ может наступить только случайно. Его можно достигнуть только при помощи более глубоко лежащих средств.

Для подготовки функции к интегрированию может быть, по-видимому, рекомендовано следующее не вполне строгое практическое правило.

Если избранное правило квадратур (5.10.1) имеет алгебраическую степень точности $m-1$, то следует стремиться к тому, чтобы интегрируемая функция f была m -кратно непрерывно дифференцируемой.

Укажем теперь некоторые правила увеличения порядка дифференцируемости.

1. Пусть интеграл имеет форму

$$\int_a^b (x-x_1)^{\alpha} \varphi(x) dx, \quad (5.10.3)$$

где x_1 есть некоторая точка, лежащая на отрезке $[a, b]$ или близко от него. Для определенности остановимся на случае, когда x_1 принадлежит

$[a, b]$. Показатель степени α будем предполагать большим -1 и не равным целому числу.

Функцию φ предположим l -кратно непрерывно дифференцируемой на $[a, b]$ и такой, что $\varphi(x_1) \neq 0$.

Когда $\alpha < 0$, интеграл будет несобственным, при $\alpha > 0$ у интегрируемой функции производные в точке x_1 не будут существовать, начиная с некоторого порядка.

Разложим $\varphi(x)$ по формуле Тейлора около точки x_1 , выделим из разложения k первых членов ($k \leq m$) и положим

$$\begin{aligned} f(x) &= (x-x_1)^\alpha \varphi(x) = f_1(x) + f_2(x), \\ f_1(x) &= (x-x_1)^\alpha \left[\varphi(x_1) + \frac{x-x_1}{1!} \varphi'(x_1) + \dots + \frac{(x-x_1)^{k-1}}{k!} \varphi^{(k-1)}(x_1) \right], \\ f_2(x) &= (x-x_1)^\alpha \left[\varphi(x) - \varphi(x_1) - \frac{(x-x_1)}{1!} \varphi'(x_1) - \dots - \frac{(x-x_1)^{k-1}}{(k-1)!} \varphi^{(k-1)}(x_1) \right], \\ \int_a^b (x-x_1)^\alpha \varphi(x) dx &= \int_a^b f_1(x) dx + \int_a^b f_2(x) dx. \end{aligned}$$

Первый из интегралов может быть легко вычислен точно. $f_2(x)$ в точке x_1 имеет порядок дифференцируемости на k единиц выше, нежели $f(x)$ и $\int_a^b f_2(x) dx$ может быть вычислен при помощи правила приближенных квадратур с лучшей точностью, чем интеграл (5.10.3).

2. Мы рассмотрели случай, когда интегрируемая функция имеет степенную особенность в одной точке. Подобные преобразования могут быть проделаны, если такие особенности будут существовать в нескольких точках $[a, b]$. Рассмотрим интеграл вида

$$\int_a^b f(x) dx = \int_a^b (x-x_1)^{\alpha_1} (x-x_2)^{\alpha_2} \dots (x-x_m)^{\alpha_m} \varphi(x) dx. \quad (5.10.4)$$

Возьмем точку x_1 , отделим соответствующий ей множитель $(x-x_1)^\alpha$ и разложим по степеням $x-x_1$ произведение остальных множителей:

$$(x-x_2)^{\alpha_2} \dots (x-x_m)^{\alpha_m} \varphi(x) = \varphi_1(x) = \varphi_1(x_1) + \frac{x-x_1}{1!} \varphi_1'(x_1) + \dots$$

Отделим в разложении k_1 первых членов и положим

$$f(x) = f_1(x) + [f(x) - f_1(x)],$$

$$f_1(x) = (x - x_1)^{\alpha_1} \left[\varphi_1(x_1) + \dots + \frac{(x - x_1)^{k_1 - 1}}{(k_1 - 1)!} \varphi_1^{(k_1 - 1)}(x_1) \right].$$

Порядок дифференцируемости в точке x_1 разности $f(x) - f_1(x)$ будет на k_1 выше, чем функции $f(x)$. Аналогично строятся разложения f в остальных точках x_j :

$$f(x) = f_j(x) + [f(x) - f_j(x)]. \quad (j = 1, 2, \dots, m).$$

После этого интеграл (5.10.4) разлагается на два:

$$\int_a^b f(x) dx = \int_a^b [f_1(x) + f_2(x) + \dots + f_m(x)] dx +$$

$$+ \int_a^b [f(x) - f_1(x) - \dots - f_m(x)] dx.$$

Из них первый вычисляется точно, во втором же слагаемом интегрируемая функция будет иметь на $[a, b]$ производные более высокого порядка, чем $f(x)$ и применение к этому интегралу квадратурной формулы должно дать более точный результат, чем для (5.10.4).

3. Степенное разложение Тейлора может быть использовано для ослабления особенностей интегрируемой функции, очевидно, всякий раз, когда интеграл имеет форму

$$\int_a^b \psi(x) \varphi(x) dx,$$

где $\psi(x)$ имеет особенность в некоторой точке отрезка $[a, b]$, и интегралы

$$\int_a^b \psi(x) (x - x_1)^j dx \quad (j = 0, 1, \dots)$$

вычисляются точно, функция же $\varphi(x)$ имеет производные достаточно высокого порядка. Таким будет, например,

$$\int_a^b (x-x_1)^\alpha \ln^p |x-x_1| \varphi(x) dx$$

($\alpha > -1$ и p есть целое число).

4. Та же идея может быть применена в том случае, когда $f(x)$ есть аналитическая функция, регулярная на $[a, b]$, но имеющая в некоторой точке x_1 , лежащей вблизи $[a, b]$, особую точку степенного типа, и мы хотим устранить эту особую точку или ослабить ее влияние на погрешность приближенного интегрирования. Вот простой пример: пусть в точке x_1 функция $f(x)$ имеет полюс порядка m и представима, следовательно, в виде

$$f(x) = \frac{\varphi(x)}{(x-x_1)^m},$$

где $\varphi(x)$ регулярна в некоторой области, содержащей в себе x_1 и $[a, b]$. Разложим $\varphi(x)$ в ряд Тейлора по степеням $x-x_1$, отделим в нем члены до степени k включительно ($k \geq m$) и положим

$$f(x) = f_1(x) + [f(x) - f_1(x)],$$

$$f_1(x) = (x-x_1)^{-m} \sum_{j=1}^k \frac{(x-x_1)^j}{j!} \varphi^{(j)}(x_1).$$

Интеграл $\int_a^b f_1(x) dx$ вычисляется точно в простых функциях. Так как для $f(x) - f_1(x)$ точка x_1 не будет особой, $\int_a^b [f - f_1] dx$ при помощи правила (5.10.1) вычислится, вообще говоря, более точно, чем $\int_a^b f dx$.

§ 5.11. СХОДИМОСТЬ КВАДРАТУРНОГО ПРОЦЕССА

5.11.1. Условия сходимости общего квадратурного процесса

Будем рассматривать квадратурный процесс, определяемый бесконечными треугольными таблицами узлов $x_k^{(n)}$ и коэффициентов $A_k^{(n)}$:

$$X = \begin{pmatrix} x_1^{(1)} & & & & \\ x_1^{(2)} & x_2^{(2)} & & & \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ x_1^{(n)} & x_2^{(n)} & \dots & x_n^{(n)} & \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}, \quad A = \begin{pmatrix} A_1^{(1)} & & & & \\ A_1^{(2)} & A_2^{(2)} & & & \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ A_1^{(n)} & A_2^{(n)} & \dots & A_n^{(n)} & \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}. \quad (5.11.1)$$

Квадратурное правило, отвечающее строкам номера n этих таблиц,

$$\int_a^b p(x)f(x)dx = \sum_{k=1}^n A_k^{(n)} f(x_k^{(n)}) + R_n(f) = Q_n(f) + R_n(f) \quad (5.11.2)$$

называется сходящимся для функции f , если

$$\lim_{n \rightarrow \infty} Q_n(f) = \lim_{n \rightarrow \infty} \sum_{k=1}^n A_k^n f(x_k^{(n)}) = \int_a^b p(x)f(x)dx.$$

Выясним сейчас, каким условиям должны удовлетворять матрицы X и A , чтобы процесс сходиллся для всех функций некоторых классов. Такую задачу мы рассмотрим в классах дифференцируемых и аналитических функций, имеющих для приложений, по-видимому, наибольший интерес. Отрезок интегрирования $[a, b]$ предполагается конечным и вес $p(x)$ — любой суммируемой на $[a, b]$ функцией, не эквивалентной нулю.

Теорема 1. Для того чтобы квадратурный процесс (5.11.2) сходиллся для всякой функции f , непрерывной на $[a, b]$, необходимо и достаточно выполнение условий:

- 1) квадратурный процесс сходится для всякого многочлена,
- 2) существует число M такое, что для всех $n=1, 2, \dots$ выполняется неравенство

$$\sum_{k=1}^n |A_k^{(n)}| \leq M. \quad (5.11.3)$$

Доказательство. В множестве непрерывных на $[a, b]$ функций может быть введена норма $\|x(t)\| = \max_t |x(t)|$, после чего такое множество станет *) полным линейным нормированным пространством или пространством банахова типа $C[a, b]$.

Интеграл

$$I(f) = \int_a^b p(x)f(x)dx$$

и квадратурная сумма

$$Q_n(f) = \sum_{k=1}^n A_k^{(n)} f(x_k^{(n)})$$

есть два линейных функционала, преобразующих пространство $C[a, b]$ в числовое пространство, которое также является пространством типа Банаха.

Для выяснения условий, при которых $Q_n(f) \rightarrow I(f)$ ($n \rightarrow \infty$), можно воспользоваться теоремой Банаха — Штейнгауза о сходимости последовательности линейных операторов.** Она говорит о том, что для сходимости $Q_n(f)$ к $I(f)$ необходимо и достаточно выполнение двух условий:

- 1) множество элементов всюду плотно в пространстве, где определены операторы;
- 2) нормы операторов ограничены в совокупности.

*) Добавление I, § 2.

**) Добавление I, § 2, теорема 2'.

За множество, всюду плотное в $C[a, b]$, может быть принято множество алгебраических многочленов, так как, по теореме Вейерштрасса, ко всякой функции f , непрерывной на $[a, b]$, можно приблизиться равномерно на $[a, b]$ и сколь угодно точно при помощи таких многочленов.

Первое условие теоремы Банаха — Штейнгауза будет выполнено в нашей задаче, если сходимость $Q_n(f) \rightarrow I(f)$ имеет место для всякого алгебраического многочлена.

Так как

$$\|Q_n(f)\| = \max_{|f| \leq 1} \left| \sum_{k=1}^n A_k^{(n)} f(x_k^{(n)}) \right| = \sum_{k=1}^n |A_k^{(n)}|,$$

второе условие об ограниченности в совокупности норм функционалов эквивалентно неравенству (5.11.3).

Этим теорема доказана.

Получим два простых следствия из нее. Рассмотрим интерполяционный квадратурный процесс. Напомним, это означает, что коэффициенты $A_k^{(n)}$ имеют следующие значения:

$$A_k^{(n)} = \int_a^b p(x) \omega_n(x) [(x - x_k^{(n)}) \omega_n'(x_k^{(n)})]^{-1} dx,$$

$$\omega_n(x) = (x - x_1^{(n)}) \dots (x - x_n^{(n)}).$$

Такой процесс определяется таблицей узлов X .

Теорема 2. Чтобы интерполяционный квадратурный процесс сходиллся для всякой непрерывной функции, необходимо и достаточно выполнение неравенства (5.11.3):

$$\sum_{k=1}^n |A_k^{(n)}| \leq M < \infty \quad (n=1, 2, \dots).$$

Доказательство. Первое условие теоремы 1 здесь, очевидно, выполняется, так как если f есть многочлен степени m , то при $n > m$ будет

$$Q_n(f) = \int_a^b p f dx.$$

Второе же условие теоремы 1 совпадает с условием теоремы 2.

Следует заметить, что хотя теорема 2 дает необходимый и достаточный признак сходимости, но этот признак нельзя, к сожалению, признать эффективным, так как весьма трудно сказать, какое заключение о расположении узлов $x_k^{(n)}$ на $[a, b]$ вытекает из неравенства (5.11.3).

Теорема 3. Если коэффициенты $A_k^{(n)}$ неотрицательны, то квадратурный процесс сходится для всякой непрерывной функции f в том и только в том случае, когда он является сходящимся для всякого алгебраического многочлена.

Доказательство. Необходимость условия является очевидной. Нужно проверить лишь его достаточность. Когда f есть многочлен нулевой степени $f \equiv 1$, должно быть

$$Q_n(1) \rightarrow \int_a^b p dx \quad (n \rightarrow \infty).$$

Поэтому множество чисел $Q_n(1)$ ограничено: $Q_n(1) \leq K$. Отсюда следует

$$\sum_{k=1}^n |A_k^{(n)}| = \sum_{k=1}^n A_k^{(n)} = Q_n(1) \leq K.$$

Таким образом, из допущения о сходимости процесса для алгебраических многочленов следует выполнение второго условия теоремы 1 и процесс будет сходиться для всякой непрерывной функции.

Отметим попутно, что теорема 7 § 5.5 о сходимости квадратурного процесса наивысшей алгебраической степени точности может быть получена как прямое следствие теоремы 3. В самом деле, если $p(x) \geq 0$, то квадратурное правило наивысшей степени точности может быть построено при всяком n . Коэффициенты $A_k^{(n)}$, как доказано в теореме 5 § 5.2, положительны.

Последовательность таких правил образует квадратурный процесс с положительными $A_k^{(n)}$. Если f есть многочлен степени m , то при $2n-1 \geq m$ будет $Q_n(f) = I(f)$ и процесс сходится для всякого многочлена. Но по теореме 3 он тогда будет сходиться для всякой непрерывной на $[a, b]$ функции.

Теперь перейдем к выяснению условий сходимости квадратурного процесса (5.11.2) на множествах дифференцируемых функций.

Будем считать, что узлы $x_k^{(n)}$ перенумерованы в порядке роста:

$$a \leq x_1^{(n)} < \dots < x_n^{(n)} \leq b,$$

и введем кусочно постоянную функцию, связанную с узлами $x_k^{(n)}$ и значениями коэффициентов $A_k^{(n)}$,

$$F_{n0}(x) = \sum_{k=1}^n A_k^{(n)} E(x - x_k^{(n)}).$$

Наряду с $F_{n0}(x)$ будем рассматривать первообразные функции для нее $F_{nr}(x)$ любого порядка r , удовлетворяющие начальным условиям $F_{nr}^{(j)}(a) = 0$ ($j = 0, 1, \dots, r-1$),

$$F_{nr}(x) = \sum_{k=1}^n A_k^{(n)} E(x - x_k^{(n)}) \frac{1}{r!} (x - x_k^{(n)})^r. \quad (5.11.4)$$

Рассмотрим множество $C_r[a, b]$ функций, имеющих непрерывные производные порядка r на $[a, b]$.

Теорема 4. Чтобы квадратурный процесс (5.11.2) сходиллся для всякой функции $f \in C_r[a, b]$ ($r \geq 1$), необходимо и достаточно выполнение условий:

1) процесс сходится для всякого многочлена;

2) существует число $M < \infty$ такое, что при всяких значениях $n = 1, 2, \dots$ выполняется условие

$$\int_a^b |F_{n, r-1}(t)| dt \leq M < \infty. \quad (5.11.5)$$

Доказательство. Сначала убедимся в необходимости условий. Необходимость первого из них очевидна и мы должны проверить лишь необходимость второго. Для этой цели воспользуемся представлением (4.8.34), характерным для функций из $C_r[a, b]$:

$$f(x) = \sum_{i=0}^{r-1} c_i (x-b)^i + (-1)^r \int_a^b g(t) E(t-x) \frac{(t-x)^{r-1}}{(r-1)!} dt, \quad (5.11.6)$$

$$c_i = \frac{1}{i!} f^{(i)}(b), \quad g(t) = f^{(r)}(t).$$

Напомним, что c_i здесь — произвольные численные параметры и $g(t)$ — любая непрерывная функция.

Остаток квадратуры

$$R_n(f) = \int_a^b p dx - \sum_{i=1}^n A_h^{(n)} f(x_h^{(n)}),$$

если в него внести вместо f ее представление (5.11.6), примет следующий вид:

$$R_n(f) = \sum_{i=0}^{r-1} c_i R_n[(x-b)^i] + (-1)^r \int_a^b g(t) \left[\int_a^t p(x) \frac{(t-x)^{r-1}}{(r-1)!} dx - F_{n, r-1}(t) \right] dt.$$

Стремление $R_n(f)$ к нулю, ввиду независимости параметров c_i и $g(t)$, равносильно выполнению следующих соотношений при $n \rightarrow \infty$:

$$R_n[(x-b)^i] \rightarrow 0 \quad (i=0, 1, \dots, r-1) \quad (5.11.7)$$

и

$$R_n^*(g) = \int_a^b g(t) \left[\int_a^t p(x) \frac{(t-x)^{r-1}}{(r-1)!} dx - F_{n, r-1}(t) \right] dt \rightarrow 0. \quad (5.11.8)$$

Остановим свое внимание на условии (5.11.8). Оно должно выполняться для всякой непрерывной функции $g(t)$. Введем на множестве таких функций норму $\|g\| = \max |g(t)|$.

$R_n^*(g)$ можно рассматривать как оператор, переводящий пространство $C[a, b]$, принадлежащее банахову типу, в одномерное числовое пространство и (5.11.8) есть условие

сходимости последовательности $R_n^*(g)$ к нулевому оператору. Здесь может быть применена теорема Банаха — Штейнгауза (см. добавление I, § 2, теорема 2), согласно которой для сходимости последовательности операторов, переводящих пространство банахова

типа в пространство того же типа, независимо от того, каким будет предельный оператор, необходимо и достаточно выполнение двух условий.

1. Последовательность операторов должна сходиться на множестве элементов, всюду плотном в пространстве, где заданы операторы.

Для нас сейчас это условие значения не имеет, так как в необходимости выполнения неравенства (5.11.5) мы убедимся, не пользуясь им.

2. Нормы операторов должны быть ограничены в совокупности.

Так как

$$\int_a^b g(t) \int_a^t p(x) \frac{(t-x)^{r-1}}{(r-1)!} dx dt$$

не зависит от n , такое требование равносильно ограниченности в совокупности норм операторов

$$\bar{R}_n(g) = \int_a^b g(t) F_{n, r-1}(t) dt.$$

Очевидно,

$$|\bar{R}_n(g)| \leq \int_a^b |F_{n, r-1}(t)| dt \max_t |g(t)| = \int_a^b |F_{n, r-1}(t)| dt \|g\|.$$

Поэтому для нормы \bar{R}_n верно неравенство

$$\|\bar{R}_n\| \leq \int_a^b |F_{n, r-1}(t)| dt.$$

Функция $F_{n, r-1}(t)$ внутри каждого из отрезков $(a, x_1^{(n)})$, $(x_1^{(n)}, x_2^{(n)})$, \dots , $(x_n^{(n)}, b)$, как видно из (5.11.4), есть многочлен и имеет либо конечное число перемен знака, либо есть тождественный нуль. Поэтому $\text{sign } F_{n, r-1}(t)$ имеет на $[a, b]$ только конечное число точек разрыва, и для любого $\varepsilon > 0$ наверное существует такая непрерывная функция $\bar{g}(t)$ ($|\bar{g}(t)| \leq 1$), что будет

$$\int_a^b |F_{n, r-1}(t)| dt - \varepsilon = \int_a^b F_{n, r-1}(t) \text{sign } F_{n, r-1}(t) dt - \varepsilon < \int_a^b F_{n, r-1}(t) \bar{g}(t) dt.$$

Для оценки $\|\bar{R}_n\|$ это даст

$$\|\bar{R}_n\| = \sup_{|g| \leq 1} \left| \int_a^b g(t) F_{n, r-1}(t) dt \right| \geq \int_a^b \bar{g}(t) F_{n, r-1}(t) dt > \int_a^b |F_{n, r-1}(t)| dt - \varepsilon,$$

и так как неравенство верно при всяких ε , то

$$\|\bar{R}_n\| \geq \int_a^b |F_{n, r-1}(t)| dt.$$

Сравнение с полученной раньше оценкой $\|\bar{R}_n\|$ сверху приводит к заключению, что

$$\|\bar{R}_n\| = \int_a^b |F_{n, r-1}(t)| dt.$$

Таким образом, второе условие теоремы и требование ограниченности в совокупности норм операторов $R_n^*(g)$ — равносильны.

Проведенные рассуждения убеждают в том, что условия теоремы являются необходимыми.

Теперь проверим достаточность условий теоремы 4. Предположим, что они выполнены. Если $g(t) = f^{(r)}(t)$ есть многочлен, то f также будет многочленом и, по первому условию, квадратурный процесс должен сходиться: $R_n(f) \rightarrow 0$. Но

$$R_n(f) = \sum_{i=0}^{r-1} c_i R_n[(x-b)^i] + (-1)^r R_n^*(g),$$

и так как

$$R_n[(x-b)^i] \rightarrow 0 \quad (n \rightarrow \infty),$$

то для любого многочлена $g(t)$ должно быть $R_n^*(g) \rightarrow 0$. Таким образом, на множестве многочленов, всюду плотном в $C[a, b]$, последовательность $R_n^*(g)$ сходится к нулю.

Условие (5.11.5) означает, что нормы операторов $\bar{R}_n(y)$ ограничены в совокупности числом M :

$$\|\bar{R}_n\| \leq M.$$

Но тогда, на что мы обращали внимание выше, будут ограничены в совокупности и нормы операторов $R_n^*(g)$:

$$\|R_n^*(g)\| \leq N.$$

Для операторов $R_n^*(g)$ выполняются оба условия видоизмененной теоремы Банаха — Штейнгауза *) и последовательность R_n^* должна сходиться к нулю на $C[a, b]$. Вместе с ней должна сходиться к нулю и последовательность операторов $R_n(f)$ при всякой функции $f \in C_r[a, b]$. Этим доказана достаточность условий теоремы 4.

Рассмотрим частный случай теоремы 4 для $r=1$ и найдем условия сходимости процесса (5.11.2) на множестве однократно непрерывно дифференцируемых функций.

$F_{n0}(t)$ есть кусочно постоянная функция со следующими значениями на интервалах между узлами:

*) Добавление 1, § 2, теорема 2'.

$$F_{n0}(t) = \begin{cases} 0 & \text{при } a \leq t < x_1^{(n)}, \\ A_1^{(n)} & \text{при } x_1^{(n)} < t < x_2^{(n)}, \\ A_1^{(n)} + A_2^{(n)} & \text{при } x_2^{(n)} < t < x_3^{(n)}, \\ \dots & \dots \\ A_1^{(n)} + \dots + A_n^{(n)} & \text{при } x_n^{(n)} < t \leq b. \end{cases}$$

Значит,

$$\begin{aligned} \int_a^b |F_{n0}(t)| dt &= |A_1^{(n)}| (x_2^{(n)} - x_1^{(n)}) + |A_1^{(n)} + A_2^{(n)}| (x_3^{(n)} - x_2^{(n)}) + \dots + \\ &+ |A_1^{(n)} + \dots + A_{n-1}^{(n)}| (x_n^{(n)} - x_{n-1}^{(n)}) + |A_1^{(n)} + \dots + A_{n-1}^{(n)} + A_n^{(n)}| (b - x_n^{(n)}). \end{aligned}$$

Отсюда следует

Теорема 5. Для сходимости квадратурного процесса (5.11.2) на множестве непрерывно дифференцируемых на $[a, b]$ функций необходимо и достаточно выполнение условий:

- 1) процесс сходится для всякого многочлена;
- 2) существует число M такое, что при $n=1, 2, \dots$ выполняются неравенства

$$|A_1^{(n)}| (x_2^{(n)} - x_1^{(n)}) + |A_1^{(n)} + A_2^{(n)}| (x_3^{(n)} - x_2^{(n)}) + \dots + |A_1^{(n)} + \dots + A_n^{(n)}| (b - x_n^{(n)}) \leq M.$$

5.11.2. Сходимость интерполяционных квадратурных процессов

Интерполяционный квадратурный процесс

$$\int_a^b p(x) f(x) dx = \sum_{k=1}^n A_k^{(n)} f(x_k^{(n)}) + R_n, \quad (5.11.9)$$

$$A_k^{(n)} = \int_a^b p(x) \frac{\omega_n(x)}{(x - x_k^{(n)}) \omega_n'(x_k^{(n)})} dx, \quad \omega_n = \prod_{j=1}^n (x - x_j^{(n)}),$$

определяется таблицей X узлов $x_k^{(n)}$ (5.11.1). Погрешность правила R_n равна интегралу от остатка интерполирования:

$$R_n = \int_a^b p(x) r_n(x) dx, \quad r_n(x) = f(x) - \sum_{k=1}^n \frac{\omega_n(x)}{(x - x_k^{(n)}) \omega_n'(x_k^{(n)})} f(x_k^{(n)}). \quad (5.11.10)$$

Многие теоремы о сходимости интерполирования могут служить источниками теорем о сходимости интерполяционных квадратурных процессов с теми же таблицами узлов $x_k^{(n)}$. В частности, если отрезок интегрирования $[a, b]$ конечный и если интерполяционный процесс на $[a, b]$ сходится равномерно, то будет сходиться и интерполяционный

квадратурный процесс (5.11.9) с той же таблицей X узлов при любой суммируемой весовой функции $p(x)$.

Приведем несколько примеров теорем такого рода. Рассмотрим правило Ньютона — Котеса

$$\int_0^1 p(x) f(x) dx \approx \sum_{k=1}^n B_k^n f\left(\frac{k}{n}\right), \quad (5.11.11)$$

$$B_k^n = \frac{(-1)^{n-k}}{nk!(n-k)!} \int_0^n \frac{t(t-1)\dots(t-n)}{t-k} dt.$$

Таблица узлов для него:

$$X = \begin{pmatrix} 0 & 1 & & & \\ 0 & \frac{1}{2} & 1 & & \\ 0 & \frac{1}{3} & \frac{2}{3} & 1 & \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \frac{1}{n} & \frac{2}{n} & \dots & 1 \end{pmatrix}. \quad (5.11.12)$$

Она, очевидно, имеет предельную функцию распределения $\mu(x) = x$. К интерполированию с таблицей узлов (5.11.12) применима теорема 2 § 4.8 и она дает возможность высказать теорему о сходимости процесса Ньютона — Котеса.

Теорема 6. Если аналитическая функция $f(z)$ регулярна в замкнутой области, ограниченной линией *)

$$\operatorname{Re} [z \ln z + (1-z) \ln(1-z)] = 1,$$

то квадратурный процесс Ньютона — Котеса для нее сходится при всякой суммируемой весовой функции $p(x)$.

Простым следствием теоремы 3 § 4.8 является

Теорема 7. Пусть отрезок интегрирования есть $[-1, 1]$. Если таблица узлов X имеет предельную функцию распределения и она является функцией Чебышева

$$\mu(x) = \pi^{-1} \int_{-1}^x \frac{dt}{\sqrt{1-t^2}},$$

то интерполяционный квадратурный процесс сходится для всякой функции $f(x)$, аналитической на отрезке $[-1, 1]$, при любом суммируемом весе $p(x)$.

Можно было бы увеличить число таких теорем, но мы ограничимся двумя приведенными примерами.

*) См. рис. 4.8.1.

§ 5.12. ВЫЧИСЛЕНИЕ НЕОПРЕДЕЛЕННОГО ИНТЕГРАЛА

5.12.1. Введение

Предположим, что на отрезке $[x_0, X]$ задана непрерывная функция $f(x)$. Всякая первообразная для нее, как известно, представима в форме

$$y(x) = y_0 + \int_{x_0}^x f(t) dt \quad (5.12.1)$$

и ее вычисление приводится к нахождению значений интеграла с переменной границей. Тот факт, что верхняя граница x является переменной и значения $y(x)$ нужно находить обычно для многих x , делает задачу неопределенного интегрирования своеобразной и побуждает для ее решения строить свои методы, учитывающие особенности этой задачи.

Пусть нужно вычислить (5.12.1) для заданной сетки значений аргумента x : $x_0 < x_1 < x_2 < \dots$. Допустим, что вычисления начаты, доведены до значения x_n и составлена приводимая в тексте таблица. Найти нужно y_{n+1} . Для этого можно воспользоваться любыми уже найденными значениями y_k ($k \leq n$), так же как любыми значениями f , которые можно применить в вычислениях.

x	y
x_0	y_0
x_1	y_1
\dots	\dots
x_n	y_n
x_{n+1}	

Поясним это немного подробнее. Когда f задана таблично, мы имеем право выбора как числа значений f , так и положения этих значений в таблице. Задача построения правила вычислений является в этом случае комбинаторной и принадлежит дискретной математике.

Если же f задана формулой, в нашем распоряжении имеются значительно большие возможности, так как мы можем произвольно избирать те точки x , в которых берутся значения f , и ограничивать себя только в числе таких значений.

Выбором значений x обычно стремятся достичь возможно высокой алгебраической степени точности вычисления y_{n+1} , что часто позволяет сделать погрешность не выше заданной границы при меньшем числе значений f , чем в других правилах нахождения y_{n+1} , и это дает возможность экономить вычислительный труд.

Все изложенное является одинаковым как в задаче определенного, так и неопределенного интегрирования, но при неопределенном интегрировании есть еще другое средство сбережения времени работы машин и труда вычислителя. Ввиду того что находить нужно многие значения $y(x)$, можно стремиться к тому, чтобы каждое значение f применялось для нахождения не одного, а нескольких значений функции $y(x)$.

Задержим внимание на этом факте и поясним его на частном правиле. Для нахождения y_{n+1} можно, например, воспользоваться равенством

$$y_{n+1} = y_n + \int_{x_n}^{x_{n+1}} f(t) dt.$$

Вычислить интеграл, стоящий справа, можно было бы при помощи почти любого из правил нахождения численного значения определенного интеграла, о которых говорилось выше. Для этого мы должны будем на отрезке интегрирования $[x_n, x_{n+1}]$ взять несколько узлов, найти значения f в них и, наконец, образовать линейную комбинацию из них с соответствующими квадратурными коэффициентами.

Это — возможный способ вычисления, и к нему нередко прибегают, но мы не будем останавливаться на такого рода методах, так как они являются простым применением к неопределенному интегрированию уже знакомых правил вычисления определенных интегралов. Эти методы имеют очевидный недостаток: значения функции f , если они не отвечают концам отрезка $[x_n, x_{n+1}]$, используются только для вычисления y_{n+1} и не участвуют в вычислениях предшествующих значений y_n, y_{n-1}, \dots , так же как и следующих значений y_{n+2}, y_{n+3}, \dots .

В дальнейшем мы будем говорить преимущественно о таких методах, которые позволяют использовать каждое значение f на нескольких шагах вычислений.

Укажем еще на возникающую в связи с неопределенным интегрированием, но имеющую, как мы увидим ниже, много более широкое значение проблему роста погрешности при вычислениях на большое число шагов. Пусть для нахождения значения y_{n+1} избрано какое-либо правило. При его применении на каждом шаге мы будем совершать некоторую погрешность. Такие погрешности постепенно будут накапливаться, и величина погрешности будет, вообще говоря, от шага к шагу увеличиваться.

Закон увеличения погрешности, если говорить в весьма общих чертах и иметь в виду не только задачу неопределенного интегрирования, но и другие аналогичные задачи, решение которых требует многократного применения одного и того же вычислительного правила, зависит как от самой задачи, так и от избранного правила вычислений. При неудачном выборе правила рост погрешности может оказаться настолько быстрым, что уже через небольшое число шагов она может стать выше допустимой границы. Более наглядным и поучительным, чем общие рассуждения, здесь, вероятно, будет приводимый ниже простой пример. Пусть на отрезке $[0, 1]$ нужно вычислить в равноотстоящих точках значения интеграла

$$y(x) = \int_0^x e^t dt = e^x - 1.$$

Для нахождения y_{n+1} воспользуемся сначала двумя предшествующими значениями функции y_n и y_{n-1} и двумя предшествующими значениями производной $y'_n = f_n$ и $y'_{n-1} = f_{n-1}$ и выполним по ним интерполирование y_{n+1} . Это есть интерполирование с двумя двукратными узлами и к нему может быть применено правило (4.7.10), которое в нашем случае будет:

$$y_{n+1} = -4y_n + 5y_{n-1} + h(4f_n + 2f_{n-1}). \quad (5.12.2)$$

Равенство является точным для всех алгебраических многочленов третьей степени. При применении оно требует знания двух начальных значений y_0 и y_1 .

Проведем вычисления с шагом $h=0,2$, считая известными $y_0=0$ и $y_1=y(0,2) \approx 0,22140$. Вычисления выполнены с запасными знаками, чтобы показать, что погрешности вызваны не недостаточным числом верных значащих цифр, а иными причинами.

x_n	y_n	погрешность
0,0	0,00000	
0,2	0,22140	
0,4	0,49152	+0,00030
0,6	1,82296	-0,00082
0,8	1,22026	+0,00528
1,0	1,74294	-0,02466

По таблице результатов видно, что погрешность быстро растет при удалении от начала таблицы. Покажем теперь, что быстрый рост вызван не большой величиной шага h и поведение погрешности не может быть улучшено при помощи уменьшения h . Уменьшим шаг вдвое, положив

$h=0,1$, и выполним вычисления вторично, приняв известными значения $y(0)=0$, $y(0,1)=0,10517$.

x_n	y_n	погрешность	x_n	y_n	погрешность
0,0	0,00000		0,6	0,81810	+0,00402
0,1	0,10517		0,7	1,03610	-0,02235
0,2	0,22139	+0,00001	0,8	1,11602	+0,10952
0,3	0,34988	-0,00002	0,9	2,01039	-0,55079
0,4	0,49165	+0,00017	1,0	-1,03251	+2,75079
0,5	0,64950	-0,00078			

Как видно, уменьшение шага в два раза позволило улучшить лишь значение $y(0,4)$, погрешности же остальных значений возросли. Рост погрешности остался столь же быстрым, как и раньше.

Наконец, чтобы убедиться в том, что рост погрешности вызван только плохими свойствами расчетного правила (5.12.2), выполним вычисление того же интеграла при помощи другого правила, которое имеет меньшую

алгебраическую степень точности, но для которого закон роста погрешности является много более благоприятным, чем для (5.12.2). Если в равенстве

$$y(x_{n+1}) = y(x_n) + \int_{x_n}^{x_{n+1}} f(t) dt$$

интеграл вычислить при помощи элементарного правила трапеций и отбросить остаточный член, получится правило

$$y_{n+1} = y_n + \frac{1}{2} h (f_n + f_{n+1}). \quad (5.12.3)$$

Оно дает верный результат в том случае, когда f есть линейная функция, и алгебраическая степень точности его ниже, чем у (5.12.2). Поэтому естественно было бы думать, что применение (5.12.3) должно для y_n дать значения менее точные, чем (5.12.2).

x_n	y_n (5.12.3)	погрешность	x_n	y_n (5.12.3)	погрешность
0,0	0,00000		0,6	0,82280	-0,00068
0,1	0,10526	-0,00009	0,7	1,01459	84
0,2	0,22159	19	0,8	1,22656	102
0,3	0,35015	29	0,9	1,46082	122
0,4	0,49223	41	1,0	1,71971	143
0,5	0,64926	54			

Для нескольких первых значений y_n это действительно так и есть, но для правила (5.12.3) погрешности растут заметно медленнее, чем для (5.12.2), и значения y_n , не близкие к началу таблицы, получаются более точными.

Приведенные примеры показывают, что вычислительные правила, даже если они имеют не низкую степень точности и могут дать малую погрешность при однократном их применении, не всегда являются пригодными при счете на большое число шагов. При выборе и построении правил вычисления необходимо обращать внимание на соответствующий им рост погрешностей, и чем большее число шагов нужно выполнить, тем более тщательным должен быть отбор правил.

Для характеристики применимости правил к многошаговым вычислениям были введены понятия «устойчивости» и «неустойчивости» правил относительно роста погрешности. Так как характер изменения погрешности, кроме вычислительного правила, зависит также от типа решаемой

задачи, то содержание этих двух понятий для каждой задачи может быть своим. Понятия еще не вполне установились, но какое бы содержание в них ни вкладывалось, «устойчивыми», по-видимому, следует считать лишь такие правила, которые позволяют найти сколь угодно точно решение задачи. Иными словами говоря, «устойчивыми» следует признать только те правила вычислений, которые позволяют построить вычислительный процесс,^{*)} сходящийся к точному решению.

Более подробно вопрос об «устойчивости» будет рассматриваться ниже при изучении методов решения дифференциальных уравнений, обыкновенных и в частных производных. Задача неопределенного интегрирования является простой в принципиальном отношении и, как будет выяснено ниже, для нее вопрос о том, какие правила вычислений следует считать «устойчивыми» и какие «неустойчивыми», выясняется достаточно просто.

Заметим также, что разделения правил вычислений лишь на два класса — «устойчивых» и «неустойчивых» относительно роста погрешностей — недостаточно, чтобы характеризовать достоинства и недостатки правил. В частности, такое разделение не дает никаких сведений о скорости роста погрешности и о потере числа верных значащих цифр.

В некоторых случаях возникает потребность говорить о более сильной или более слабой неустойчивости и бывает целесообразно ввести численную меру неустойчивости, состоящую из одного числа или нескольких чисел. Аналогичные вопросы возникают и при сравнении между собой «устойчивых» правил вычислений.

5.12.2. Погрешность вычислений и сходимость

Всюду ниже мы будем считать, что неопределенный интеграл (5.12.1) нужно найти для равноотстоящих значений аргумента $x_k = x_0 + kh$ ($h > 0$).

Для пояснения вопроса о росте погрешности при неопределенном интегрировании достаточно рассмотреть вычислительное правило

$$y_{n+1} = \sum_{i=0}^p A_i y_{n-i} + h \sum_{j=1}^m B_{nj} f(\xi_{nj}). \quad (5.12.4)$$

Оно должно быть дополнено начальными значениями y_0, y_1, \dots, y_p функции $y(x)$ и, кроме того, — указано правило округления при вычислении правой части. Если операцию округления обозначить фигурными скобками, то действительное расчетное правило будет

^{*)} Равенство (5.12.2) вместе с указанием способа нахождения начальных значений y_0 и y_1 и указанием закона округления при вычислении правой части образует вычислительный алгоритм. Последовательность же таких алгоритмов, отвечающих последовательности значений h , сходящейся к нулю, составляет вычислительный процесс.

$$y_{n+1} = \left\{ \sum_{i=0}^p A_i y_{n-i} + h \sum_{j=1}^m B_{nj} f(\xi_{nj}) \right\}_n. \quad (5.12.5)$$

Введем понятие о погрешности формулы (5.12.4). Если в нее вместо y_k подставить точные значения $y(x_k)$ первообразной, равенство не будет выполняться и в правую часть его необходимо ввести дополнительный член

$$y(x_{n+1}) = \sum_{i=0}^p A_i y(x_{n-i}) + h \sum_{j=1}^m B_{nj} f(\xi_{nj}) + r_n. \quad (5.12.6)$$

Величина r_n называется погрешностью формулы (5.12.4). Рассмотрим теперь погрешность вычисления

$$\varepsilon_n = y(x_n) - y_n$$

и получим уравнение для нее. Приближенное значение y_n находится по правилу (5.12.5), которое, если погрешность округления обозначить $-\alpha_n$, равносильно правилу

$$y_{n+1} = \sum_{i=0}^p A_i y_{n-i} + h \sum_{j=1}^m B_{nj} f(\xi_{nj}) - \alpha_n. \quad (5.12.7)$$

Вычитание его из (5.12.6) дает

$$\varepsilon_{n+1} = \sum_{i=0}^p A_i \varepsilon_{n-i} + r_n + \alpha_n. \quad (5.12.8)$$

Это есть линейное неоднородное уравнение в конечных разностях порядка $p+1$ с постоянными коэффициентами. При его решении значения погрешностей $\varepsilon_0, \varepsilon_1, \dots, \varepsilon_p$, отвечающие приближенным значениям y_k ($k=0, 1, \dots, p$), образующим начало расчетной таблицы, мы должны предполагать известными. Все следующие значения ε_n ($n > p$) должны быть найдены из (5.12.8). Положив там $n=p$, мы найдем ε_{p+1} в форме линейной комбинации ε_k ($k=0, 1, \dots, p$) и величины $r_p + \alpha_p$. Пользуясь этим результатом и положив в (5.12.8) $n=p+1$, найдем ε_{p+2} в виде линейной комбинации тех же начальных погрешностей ε_k ($k=0, 1, \dots, p$) и величин $r_p + \alpha_p, r_{p+1} + \alpha_{p+1}$ и т. д.

При помощи уравнения (5.12.8) погрешность ε_n будет найдена как линейная однородная функция начальных погрешностей ε_k ($k \leq p$) и величин $r_p + \alpha_p, \dots, r_{n-1} + \alpha_{n-1}$. Коэффициенты этой функции будут, очевидно, зависеть от n :

$$\varepsilon_n = \Gamma_n^0 \varepsilon_0 + \Gamma_n^1 \varepsilon_1 + \dots + \Gamma_n^p \varepsilon_p + \sum_{j=p}^{n-1} G_n^j (r_j + \alpha_j). \quad (5.12.9)$$

В написанном равенстве можно считать $n=0, 1, 2, \dots$, если условиться сумму $\sum_{j=p}^{n-1}$ заменять нулем в тех случаях, когда верхняя граница суммирования $n-1$ становится меньше p .

Γ_n^i называется функцией влияния или гриновой функцией начального значения ϵ_i . Аналогично, G_n^j называют функцией влияния значения $r_j + \alpha_j$ свободного члена уравнения. Γ_n^i и G_n^j являются частными решениями уравнения вида (5.12.8) с указываемыми ниже начальными условиями.

Возьмем сначала Γ_n^0 . Положим свободный член уравнения тождественно равным нулю:

$$r_n + \alpha_n = 0 \quad (n=p, p+1, \dots)$$

и рассмотрим однородное уравнение

$$L(\epsilon_n) = \epsilon_{n+1} - \sum_{i=0}^p A_i \epsilon_{n-i} = 0.$$

Затем будем считать $\epsilon_0=1, \epsilon_1=\dots=\epsilon_p=0$. Тогда, как видно из (5.12.9), будет $\epsilon_n = \Gamma_n^0$. Значит, можно сказать, что Γ_n^0 есть решение однородного уравнения $L(\epsilon_n)=0$ с начальными условиями

$$\epsilon_0=1, \quad \epsilon_1=\epsilon_2=\dots=\epsilon_p=0.$$

Функция влияния Γ_n^1 начального значения ϵ_1 есть решение однородного уравнения $L(\epsilon_n)=0$ с начальными значениями

$$\epsilon_0=0, \quad \epsilon_1=1, \quad \epsilon_2=\dots=\epsilon_p=0, \dots$$

Для нас особый интерес будет иметь функция влияния Γ_n^p , принадлежащая последнему начальному значению ϵ_p . Она есть решение следующей задачи:

$$L(\epsilon_n)=0, \quad \epsilon_0=\epsilon_1=\dots=\epsilon_{p-1}=0, \quad \epsilon_p=1. \quad (5.12.10)$$

В некоторых вопросах удобно продолжить Γ_n^p на отрицательные значения n , считая $\Gamma_n^p=0$ при $n=-1, -2, \dots$. Через несколько строк будет показано, что Γ_n^p для разностного уравнения с постоянными коэффициентами, каким является (5.12.8), тесно связана с функцией влияния G_n^i , к рассмотрению которой мы переходим. Если в (5.12.9) положить

$$\epsilon_0=\dots=\epsilon_p, \quad r_j + \alpha_j = \begin{cases} 0 & \text{при } j \neq i, \\ 1 & \text{при } j = i, \end{cases}$$

то получится $\epsilon_n = G_n^i$, и поэтому G_n^i есть решение следующей задачи:

$$L(\varepsilon_n) = \delta_n^i \quad (n, i \geq p), \quad \varepsilon_0 = \dots = \varepsilon_p = 0. \quad (5.12.11)$$

Здесь δ_n^i — символ Кронекера. При $n \leq i-1$ правая часть $\delta_n^i = 0$, уравнение будет однородным $L(\varepsilon_n) = 0$ и, так как начальные значения $\varepsilon_k = 0$ ($k=0, 1, \dots, p$), будет $\varepsilon_n = 0$ при $n < i$. При $n = i$ уравнение (5.12.11) даст

$$L(\varepsilon_n)|_{n=i} = \varepsilon_{i+1} - \sum_{j=0}^p A_j \varepsilon_{i-j} = \varepsilon_{i+1} = 1.$$

Для $n > i$ уравнение (5.12.11) вновь становится однородным, и G_n^i будет решением задачи

$$L(\varepsilon_n) = 0 \quad (n > i), \quad \varepsilon_{i-p+1} = 0, \varepsilon_{i-p+2} = 0, \dots, \varepsilon_i = 0, \varepsilon_{i+1} = 1.$$

Такая задача отличается от задачи (5.12.10), определяющей Γ_n^p , лишь сдвигом по оси n на $i-p+1$ единиц и поэтому $G_n^i = \Gamma_{n+p-i-1}^p$.

Это дает возможность записать ε_n в следующей более удобной для нас форме:

$$\varepsilon_n = \Gamma_n^0 \varepsilon_0 + \Gamma_n^1 \varepsilon_1 + \dots + \Gamma_n^p \varepsilon_p + \sum_{j=p}^{n-1} \Gamma_{n+p-j-1}^p (r_j + \alpha_j). \quad (5.12.12)$$

Сообразно с представлениями (5.12.9) и (5.12.12) разложим погрешность на три части:

$$\varepsilon_n = E_n + E_n' + E_n'', \quad (5.12.13)$$

$$E_n = \sum_{i=0}^p \Gamma_n^i \varepsilon_i, \quad E_n' = \sum_{j=p}^{n-1} \Gamma_{n+p-j-1}^p \alpha_j, \quad E_n'' = \sum_{j=p}^{n-1} \Gamma_{n+p-j-1}^p r_j.$$

Первая из них E_n дает часть погрешности ε_n , происходящую от неточности значений y_0, y_1, \dots, y_p , составляющих начало расчетной таблицы. E_n является решением однородного уравнения $L(E_n) = 0$ с начальными значениями $E_n = \varepsilon_n$ ($n=0, 1, \dots, p$).

Второе слагаемое E_n' учитывает влияние на ε_n погрешностей округлений α_j и есть решение неоднородного уравнения $L(E_n') = \alpha_n$ с нулевыми начальными значениями $E_n' = 0$ ($n=0, 1, \dots, p$).

Наконец, E_n'' дает часть погрешности ε_n , происходящую от неточностей r_n расчетной формулы. E_n'' есть решение неоднородного уравнения $L(E_n'') = r_n$ с нулевыми начальными значениями $E_n'' = 0$ ($n=0, 1, \dots, p$).

Ниже будут рассматриваться условия равномерной сходимости при $h \rightarrow 0$ приближенно вычисленного неопределенного интеграла с таблицей значений y_n ($n=0, 1, \dots, N$) к точным значениям $y(x_n)$ этой функции.

За расстояние между ними должна быть принята величина

$$\rho(y, y_n) = \max_n |y(x_n) - y_n|,$$

и нам нужно выяснить условия, при которых будет $\rho(y, y_n) \rightarrow 0$ ($h \rightarrow 0$). Ввиду независимости величин ϵ_i ($i=0, 1, \dots, p$), r_n и α_n каждая из трех частей E_n, E_n', E_n'' должна в принятой метрике стремиться к нулю:

$$\begin{aligned} \max_n |E_n| \rightarrow 0, \quad \max_n |E_n'| \rightarrow 0, \quad \max_n |E_n''| \rightarrow 0 \quad (5.12.14) \\ (h \rightarrow 0). \end{aligned}$$

Величины E_n и E_n' зависят от погрешностей ϵ_i ($i=0, 1, \dots, p$) и α_n ($n > p$), и совершенно ясно, что при всяком фиксированном h точность вычислений начальных значений y_0, y_1, \dots, y_p и правых частей (5.12.5) можно увеличить настолько, чтобы сделать $\max_n |E_n|$ и $\max_n |E_n'|$ сколь

угодно малыми. Поэтому за счет увеличения точности вычисления, принципиально говоря, всегда можно добиться того, чтобы при $h \rightarrow 0$ было $\max_n |E_n| \rightarrow 0$ и $\max_n |E_n'| \rightarrow 0$. Вопрос здесь заключается лишь в том, на-

сколько быстро для этого необходимо увеличивать точность при убывании h . Если при $h \rightarrow 0$ погрешности ϵ_i ($i \leq p$) и α_n должны быстро убывать, то такая вычислительная схема будет малоприменимой для приложений.

В этом отношении нужно отдать предпочтение тем вычислительным схемам, для которых точность вычисления y_i ($i \leq p$) и правых частей (5.12.5) должна возрастать возможно медленнее при $h \rightarrow 0$, если мы хотим, чтобы было $\max_n |E_n| \rightarrow 0$ и $\max_n |E_n'| \rightarrow 0$.

Рассмотрим этот вопрос более подробно и начнем с изучения E_n . Если считать, что начальные погрешности ϵ_i ($i \leq p$) ограничены по абсолютной величине числом ϵ и в остальном произвольны, для E_n из (5.12.13) вытекает следующая точная оценка:

$$|E_n| \leq \epsilon \sum_{i=0}^p |\Gamma_n^i|.$$

Будем считать, что правило (5.12.4) является точным для того случая, когда y есть постоянная величина *) и $f \equiv 0$. Коэффициенты A_k такого правила должны выполнять равенство

*) Если правило вычислений не обладает этим свойством, то это правило, по-видимому, не имеет значения.

$$\sum_{i=0}^p A_i = 1.$$

Это означает, что $\varepsilon_n = 1$ есть решение однородного уравнения $L(\varepsilon_n) = 0$. Для него представление (5.12.12) примет форму $1 = \Gamma_n^0 + \Gamma_n^1 + \dots + \Gamma_n^p$. Отсюда следует, что при любых значениях n выполняется неравенство

$$\sum_{i=0}^p |\Gamma_n^i| \geq 1.$$

В отношении порядка роста E_n при неограниченном увеличении n наиболее благоприятным является случай, когда сумма $\sum_{i=0}^p |\Gamma_n^i|$ будет ограниченной.

В связи с изложенным выше в задаче вычисления неопределенного интеграла целесообразно, по-видимому, принять следующее определение.

Правило вычислений (5.12.4) называется устойчивым относительно погрешностей начальных значений ε_i ($i \leq p$), если существует число M такое, что при любых значениях $n \geq p$ будет выполняться неравенство

$$|E_n| = \left| \sum_{i=0}^p \Gamma_n^i \varepsilon_i \right| \leq M\varepsilon,$$

если только $|\varepsilon_i| \leq \varepsilon$ ($i \leq p$).

Можно легко установить признак устойчивости. Общее решение *) однородного линейного уравнения с постоянными коэффициентами $L(\varepsilon_n) = \varepsilon_{n+1} - \sum_{i=0}^p A_i \varepsilon_{n-i} = 0$ определяется алгебраическим уравнением степени $p+1$

$$\lambda^{p+1} - \sum_{i=0}^p A_i \lambda^{p-i} = \lambda^{p+1} - A_0 \lambda^p - A_1 \lambda^{p-1} - \dots - A_p = 0. \quad (5.12.15)$$

Если корни уравнения есть $\lambda_1, \lambda_2, \dots, \lambda_m$ и кратности их равны соответственно k_1, k_2, \dots, k_m , тогда функции $\lambda_i^n n_j$ ($j=0, 1, \dots, k_i-1$; $i=1, 2, \dots, m$) образуют фундаментальную систему решений однородного уравнения $L(\varepsilon_n) = 0$. Всякое решение уравнения является их линейной комбинацией.

*) Необходимые сведения из теории уравнений с конечными разностями приведены в добавлении IV.

С другой стороны, функции влияния Γ_n^i ($i=0, 1, \dots, p$) начальных значений также образуют, очевидно, фундаментальную систему и получаются из решений $\lambda_i^n n^j$ линейным преобразованием с неособенной матрицей.

Ограниченность суммы $\sum_{i=0}^n |\Gamma_n^i|$ равносильна ограниченности функций Γ_n^i ($i=0, 1, \dots, p$) и, следовательно, равносильна ограниченности при $n=1, 2, \dots$ решений $\lambda_i^n n^j$ ($j=0, 1, \dots, k_i-1$; $i=1, 2, \dots, m$), что возможно в том и только в том случае, когда среди λ_i нет чисел, больших по модулю единицы, и если $|\lambda_i|=1$, то тогда $k_i=1$. Отсюда следует теорема, дающая нужный признак.

Теорема 1. Для того чтобы правило (5.12.4) было устойчивым относительно погрешностей начальных значений y_i ($i=0, 1, \dots, p$), необходимо и достаточно выполнение условий:

1) среди корней λ_i уравнения

$$\lambda^{p+1} = \sum_{i=0}^p A_i \lambda^{p-i}$$

нет больших единицы по модулю;

2) корни уравнения, равные по модулю единице, являются простыми.

Перейдем теперь к части погрешности, происходящей от ошибок округлений α_n :

$$E_n' = \sum_{j=p}^{n-1} \Gamma_{n+p-j-1}^p \alpha_j. \quad (5.12.16)$$

Предположим, что верхняя грань α для погрешности округлений указана одинаковой для всех шагов вычислений: $|\alpha_n| \leq \alpha$ ($n \geq p$).

$$|E_n'| \leq \alpha \sum_{j=p}^{n-1} |\Gamma_{n+p-j-1}^p|,$$

$$\max_n |E_n'| \leq \alpha \sum_{j=p}^{N-1} |\Gamma_{N+p-j-1}^p| = \alpha \sum_{j=p}^{N-1} |\Gamma_j^p|. \quad (5.12.17)$$

Когда $h \rightarrow 0$, N будет неограниченно возрастать. Значение суммы $\sum_{j=p}^{N-1} |\Gamma_j^p|$ будет зависеть от поведения Γ_j^p при неограниченном росте j .

Рассмотрим $p+1$ функций

$$\Gamma_n^p, \Gamma_{n+1}^p, \dots, \Gamma_{n+p}^p. \quad (5.12.18)$$

Они являются решениями однородного уравнения $L(\epsilon_n)=0$. При $n=0, 1, \dots, p$ их значения образуют таблицу

$$\begin{pmatrix} 0 & 0 & \dots & 0 & 0 & 1 \\ 0 & 0 & \dots & 0 & 1 & \Gamma_{p+1}^p \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & \Gamma_{p+1}^p & \dots & \Gamma_{2p-2}^p & \Gamma_{2p-1}^p & \Gamma_{2p}^p \end{pmatrix},$$

определитель которой отличен от нуля. Поэтому решения (5.12.18) составляют фундаментальную систему и связаны с $\Gamma_n^0, \Gamma_n^1, \dots, \Gamma_n^p$, а следовательно, и с решениями $\lambda_i^n n^j$ неособенными линейными преобразованиями. Их ограниченность при $n \rightarrow \infty$ равносильна ограниченности Γ_n^i ($i \leq p$) и $\lambda_i^n n^j$ ($j=0, 1, \dots, k_i-1; i=1, 2, \dots, m$).

Проведенные сейчас и выше рассуждения позволяют утверждать, что для возможно медленного роста при $N \rightarrow \infty$ суммы $\sum_{k=p}^{N-1} |\Gamma_k^p|$ наиболее благоприятным будет случай, когда все Γ_k^p , а значит и $\lambda_i^n n^j$, будут ограничены. При этом условии сумма $\sum_{k=p}^{N-1} |\Gamma_k^p|$ будет величиной порядка $O(N)$.

Изложенные соображения побуждают принять приводимое ниже определение.

Правило вычислений (5.12.4) называется устойчивым относительно погрешностей округлений α_n , если существует число M_1 , не зависящее от h , такое, что при всяких $N > p$ будет выполнено неравенство

$$|E_n'| \leq M_1 N \alpha \quad (n=p, p+1, \dots, N-1). \quad (5.12.19)$$

Просто доказывается приводимая ниже теорема, дающая признак такой устойчивости.

Теорема 2. Для того чтобы правило (5.12.4) было устойчиво относительно погрешности округлений α_n , достаточно выполнение условий:

1) уравнение $\lambda^{p+1} - \sum_{i=0}^p A_i \lambda^{p-i} = 0$ не имеет корней по модулю больших единицы и

2) корни, равные единице по модулю, являются простыми.

Доказательство. Действительно, при выполнении условий теоремы будут ограничены решения $\lambda_i^n n^j$ ($j=0, 1, \dots, k_i-1$; $i=1, 2, \dots, m$), а вместе с ними будут ограничены и функции влияния Γ_n^i ($i \leq p$), в частности, Γ_n^p :

$$|\Gamma_n^p| \leq M_1.$$

Из (5.12.17) и отсюда следует утверждение теоремы.

Рассмотрим, наконец, третью часть погрешности, которая зависит от ошибок r_n самого вычислительного правила (5.12.4):

$$E_n'' = \sum_{j=p}^{n-1} \Gamma_{n+p-j-1}^p r_j.$$

Напомним, что при помощи увеличения точности вычисления начальных значений y_i ($i \leq p$) и уменьшения погрешностей округления α_n можно всегда добиться того, чтобы при $h \rightarrow 0$ было

$$\max_n |E_n| \rightarrow 0 \quad \text{и} \quad \max_n |E_n'| \rightarrow 0.$$

Это дает нам право сказать, что правило (5.12.4) наверное допускает равномерно сходящийся вычислительный процесс, если

$$\max_n |E_n''| \rightarrow 0 \quad (h \rightarrow 0). \quad (5.12.20)$$

Предположим, что $r=r(h)$ есть верхняя граница для абсолютной величины погрешности r_n при всяких n ($p \leq n \leq N-1$):

$$|r_n| \leq r=r(h).$$

Для E_n'' верны оценки

$$|E_n''| \leq r \sum_{i=p}^{n-1} |\Gamma_{n+p-j-1}^p| = r \sum_{k=p}^{n-1} |\Gamma_k^p|$$

и

$$\max_n |E_n''| \leq r \sum_{k=p}^{N-1} |\Gamma_k^p|. \quad (5.12.21)$$

Отсюда, очевидно, следует

Теорема 3. Если при $h \rightarrow 0$ будет

$$r(h) \sum_{k=p}^{N-1} |\Gamma_k^p| \rightarrow 0,$$

то правило (5.12.4) допускает равномерно сходящийся вычислительный процесс.

Рассмотрим еще частный случай этой теоремы.

Пусть правило (5.12.4) таково, что для соответствующего ему уравнения $\lambda^{p+1} = \sum_{i=0}^p A_i \lambda^{p-i}$ выполняются условия теорем 1 и 2. Тогда, как выяснилось выше, существует такое число M_1 , что при всяких $k \geq p$ будет

$$|\Gamma_k^p| \leq M_1.$$

Отсюда и из (5.12.21) получается оценка

$$\max_n |E_n''| \leq M_1(N-p)r \leq M_1 r N \leq M_1 \frac{r}{h} (X-x_0),$$

из которой вытекает

Теорема 4. Если уравнение $\lambda^{p+1} = \sum_{i=0}^p A_i \lambda^{p-i}$ не имеет корней, модули которых больше 1, и корни его, по модулю равные единице, являются однократными, то правило (5.12.4) допускает равномерно сходящийся вычислительный процесс всякий раз, когда

$$\frac{r(h)}{h} \rightarrow 0 \quad \text{при } h \rightarrow 0.$$

Приведем еще одну простую теорему, дающую достаточное условие устойчивости расчетного правила (5.12.5).

Теорема 5. Если коэффициенты правила (5.12.5) неотрицательны: $A_i \geq 0$ ($i=0, 1, \dots, p$) и $\sum_{i=0}^p A_i = 1$, то правило устойчиво относительно роста погрешности.

Возьмем соответствующее правилу однородное разностное уравнение

$$y_{n+1} = \sum_{i=0}^p A_i y_{n-i}.$$

Отсюда получается следующая оценка:

$$|y_{n+1}| \leq \sum_{i=0}^p A_i |y_{n-i}| \leq \max_{n-p \leq k \leq n} |y_k| \sum_{i=0}^p A_i = \max_{n-p \leq k \leq n} |y_k|.$$

Если применить оценку $|y_{n+1}| \leq \max_{n-p \leq k \leq n} |y_k|$ несколько раз, то можно прийти к заключению, что при любом n должно быть $|y_n| \leq \max_{0 \leq k \leq p} |y_k|$, иначе говоря, все значения решения однородного уравнения не больше максимального из модулей начальных значений y_0, y_1, \dots, y_p . Отсюда следует ограниченность абсолютных значений всех функций влияния начальных значений единицей:

$$|\Gamma_n^i| \leq 1 \quad (i=0, 1, \dots, p; n=0, 1, 2, \dots).$$

Это, как было выяснено выше, позволяет утверждать, что среди корней уравнения

$$\lambda^{p+1} = \sum_{i=0}^p A_i \lambda^{p-i}$$

нет по модулю больших единицы и корни, равные единице по модулю, являются простыми и, стало быть, правило (5.12.5) устойчиво.

§ 5.13. ПОНЯТИЕ О НЕКОТОРЫХ ЧАСТНЫХ МЕТОДАХ ВЫЧИСЛЕНИЯ НЕОПРЕДЕЛЕННОГО ИНТЕГРАЛА

5.13.1. Интегрирование функции, заданной таблицей значений

Пусть на отрезке $[x_0, X]$ в равноотстоящих точках $x_n = x_0 + nh$ ($n=0, 1, \dots, N$; $x_0 + Nh \leq X < x_0 + (N+1)h$) известны значения функции f и по ним нужно вычислить значения интеграла

$$y(x) = y_0 + \int_{x_0}^x f(t) dt \quad (5.13.1)$$

в тех же равноотстоящих точках $x_n = x_0 + nh$, где задана функция f .

Рассмотрим сначала задачу продолжения уже начатой таблицы. Вопросы о составлении начала и конца таблицы будут изучаться позже. Предположим, что вычисления доведены до узла $x_n = x_0 + nh$ и $y(x_n)$ есть последнее найденное значение функции $y(x)$. Для нахождения следующего значения $y(x_{n+1})$ мы могли бы воспользоваться любыми известными значениями $y(x_k)$ ($k \leq n$) и какими угодно табличными значениями f . Мы остановимся на методах, когда для нахождения $y(x_{n+1})$ пользуются лишь одним предыдущим значением $y(x_n)$. Точное представление $y(x_{n+1})$ через $y(x_n)$ и функцию f есть

$$y(x_{n+1}) = y(x_n) + \int_{x_n}^{x_{n+1}} f(t) dt.$$

Чтобы воспользоваться им, нужно знать $f(t)$ всюду на отрезке $[x_n, x_{n+1}]$. Точные значения $f(t)$ нам не известны, но мы можем найти $f(t)$ приближенно, выполнив интерполирование ее на $[x_n, x_{n+1}]$ по заданной таблице значений f . Для интерполирования целесообразно привлечь табличные узлы, ближайшие к отрезку $[x_n, x_{n+1}]$, и взять одинаковое число их с каждой стороны от отрезка. В качестве аппарата представления интерполирующего многочлена может быть использована формула Ньютона — Бесселя (4.4.6). В нашем случае эта формула может быть записана, если считать f непрерывно дифференцируемой $2k+2$ раза на $[x_n - kh, x_n + (k+1)h]$, в виде:

$$\begin{aligned} f(t) = f(x_n + uh) = & \frac{f_n + f_{n+1}}{2} + \frac{u-0,5}{1!} \Delta f_n + \frac{u(u-1)}{2!} \cdot \frac{\Delta^2 f_{n-1} + \Delta^2 f_n}{2} + \\ & + \frac{(u-0,5)u(u-1)}{3!} \Delta^3 f_{n-1} + \dots + \frac{(u+k-1) \dots (u-k)}{(2k)!} \times \\ & \times \frac{\Delta^{2k} f_{n-k} + \Delta^{2k} f_{n-k+1}}{2} + \frac{(u-0,5)(u+k-1) \dots (u-k)}{(2k+1)!} \Delta^{2k+1} f_{n-k} + r(t), \\ r(t) = & h^{2k+2} \frac{(u+k)(u+k-1) \dots (u-k-1)}{(2k+2)!} f^{(2k+2)}(\xi), \end{aligned}$$

$$x_n - kh < \xi < x_n + (k+1)h.$$

Подстановка такого представления f в интеграл

$$\int_{x_n}^{x_{n+1}} f(t) dt = h \int_0^1 f(x_n + uh) du$$

и несложные вычисления дадут следующее выражение для $y(x_{n+1})$:

$$\begin{aligned} y(x_{n+1}) = y(x_n) + h \left[\frac{f_n + f_{n+1}}{2} - \frac{1}{12} \cdot \frac{\Delta^2 f_{n-1} + \Delta^2 f_n}{2} + \right. \\ + \frac{11}{720} \cdot \frac{\Delta^4 f_{n-2} + \Delta^4 f_{n-1}}{2} - \frac{191}{60480} \cdot \frac{\Delta^6 f_{n-3} + \Delta^6 f_{n-2}}{2} + \\ \left. + \dots + C_k \frac{\Delta^{2k} f_{n-k} + \Delta^{2k} f_{n-k+1}}{2} \right] + R_{n,k}, \end{aligned} \quad (5.13.2)$$

$$\begin{aligned}
C_k &= \frac{1}{(2k)!} \int_0^1 (u+k-1) \dots (u-k) du, \\
R_{n,k} &= h \int_0^1 r(x_n + uh) du = \\
&= \frac{h^{2k+3}}{(2k+2)!} \int_0^1 (u+k)(u+k-1) \dots (u-k-1) f^{(2k+2)}(\xi) du = \\
&= \frac{h^{2k+3}}{(2k+2)!} f^{(2k+2)}(\eta) \int_0^1 (u+k)(u+k-1) \dots (u-k-1) du, \\
x_n - kh &< \eta < x_n + (k+1)h.
\end{aligned}$$

При переходе к последней части равенства было вынесено за знак интеграла среднее значение производной $f^{(2k+2)}$, что можно сделать, так как ядро интеграла $(u+k) \dots (u-k-1)$ сохраняет знак на отрезке $[0, 1]$.

Если в (5.13.2) отбросить неизвестный остаток $R_{n,k}$, получится приближенное расчетное правило. Когда мы хотим применить его с самого начала вычислений для нахождения значений y_1, y_2, \dots, y_k , мы должны выйти налево из интервала $[x_0, X]$ и найти дополнительно значения $f_{-1} = f(x_0 - h), f_{-2} = f(x_0 - 2h), \dots, f_{-k} = f(x_0 - kh)$.

Если мы не имеем такой возможности или хотим избежать затраты труда на дополнительные вычисления, мы должны будем при построении правила вычисления изменить в интерполировании выбор узлов. Так, например, имея целью построить правило для нахождения $y(x_1)$:

$$y(x_1) = y(x_0) + \int_{x_0}^{x_1} f(t) dt,$$

можно для интерполирования f на $[x_0, x_1]$ воспользоваться правилом Ньютона для интерполирования в начале таблицы (4.4.1):

$$\begin{aligned}
f(t) = f(x_0 + uh) &= f_0 + \frac{u}{1!} \Delta f_0 + \frac{u(u-1)}{2!} \Delta^2 f_0 + \\
&+ \frac{u(u-1)(u-2)}{3!} \Delta^3 f_0 + \dots + \frac{u(u-1) \dots (u-k+1)}{k!} \Delta^k f_0 + r,
\end{aligned}$$

$$r = h^{k+1} \frac{u(u-1)\dots(u-k)}{(k+1)!} f^{(k+1)}(\xi).$$

Подстановка в интеграл приведет к равенству

$$\begin{aligned} y(x_1) = y(x_0) + h \left[\frac{f_0 + f_1}{2} - \frac{1}{12} \Delta^2 f_0 + \frac{1}{24} \Delta^3 f_0 - \frac{19}{720} \Delta^4 f_0 + \frac{1}{160} \Delta^5 f_0 + \right. \\ \left. + \dots + C_k \Delta^k f_0 \right] + R_{n,k} \quad (k \geq 2), \\ C_k = \frac{1}{k!} \int_0^1 u(u-1)\dots(u-k+1) du, \\ R_{n,k} = C_{k+1} h^{k+2} f^{(k+1)}(\xi), \quad x_0 < \xi < x_k. \end{aligned} \quad (5.13.3)$$

После отбрасывания остатка $R_{n,k}$ получится приближенное равенство, позволяющее вычислить $y(x_1)$ по $y(x_0)$ и значениям f в точках, не выходящих налево за x_0 .

То же равенство с заменой x_0 на x_1 позволит вычислить $y(x_2)$ и т. д. Аналогичное можно сказать о вычислениях вблизи конечной точки x_N . При помощи правила (5.13.2) можно найти $y(x_n)$ до $y(x_{N-k})$ включительно. Для вычисления значений $y(x_{N-k+1}), \dots, y(x_N)$ при помощи того же правила необходимо было бы вычислить значения f_{N+1}, \dots, f_{N+k} . Чтобы избежать этого, можно воспользоваться правилом Ньютона для интерполирования в конце таблицы (4.4.3). Вычисления, сходные с сделанными выше для точки x_0 , приведут к следующему результату:

$$\begin{aligned} y(x_N) = y(x_{N-1}) + h \left[\frac{f_N + f_{N-1}}{2} - \frac{1}{12} \Delta^2 f_{N-2} - \frac{1}{24} \Delta^3 f_{N-3} - \frac{19}{720} \Delta^4 f_{N-4} - \right. \\ \left. - \frac{1}{160} \Delta^5 f_{N-5} - \dots - (-1)^{k-1} C_k \Delta^k f_{N-k} \right] + R_{n,k} \quad (k \geq 2). \end{aligned}$$

Это правило может быть применено к нахождению

$$y(x_{N-k+1}), \dots, y(x_N).$$

5.13.2. Вычисление при помощи периодически расположенных узлов

При вычислении интеграла (5.13.1) наибольшего количества труда требует обычно вычисление значений функции f . Так как находить приходится чаще всего много значений $y(x)$, можно значительно сэкономить в работе, если каждое значение применять для нахождения не одного, а многих значений $y(x)$. Этого можно добиться, если узлы, в которых вычисляются f , расположить на оси x периодически с периодом h и на каждом шаге вычислений брать значения f в сходственных точках нескольких промежутков. Поясним эту мысль более подробно.

Остановимся опять на том случае, когда для вычисления $y(x_{n+1})$ берется только одно предыдущее значение функции.

$$y(x_{n+1}) = y(x_n) + \int_{x_n}^{x_{n+1}} f(t) dt.$$

Вычислению здесь подлежит интеграл

$$\int_{x_n}^{x_{n+1}} f(t) dt.$$

Предположим, что для нахождения его значений на отрезке $[x_n, x_{n+1})$ взято m узлов $\alpha, \beta, \dots, \lambda$: $x_n \leq \alpha < \beta < \dots < \lambda < x_{n+1}$. Эти узлы назовем основными. Дополнительно возьмем еще

$$\begin{array}{ll} a \text{ узлов } \alpha + p_i h & (i=1, 2, \dots, a), \text{ сходственных } \alpha, \\ b \text{ узлов } \beta + q_i h & (i=1, 2, \dots, b), \text{ сходственных } \beta, \\ \dots & \dots \\ l \text{ узлов } \lambda + t_i h & (i=1, 2, \dots, l), \text{ сходственных } \lambda. \end{array}$$

Положение дополнительных узлов в интервалах $[x_k, x_{k+1})$ определяется числами p_i, q_i, \dots, t_i . Они могут иметь какие угодно целые значения, отличные от нуля. Их мы будем считать избранными и фиксированными. Обозначим $N+1$ общее число узлов: $m+a+b+\dots+l=N+1$ и рассмотрим правило вычислений

$$\int_{x_n}^{x_{n+1}} f(t) dt \approx A_0 f(\alpha) + \sum_{i=1}^a A_i f(\alpha + p_i h) + \dots + L_0 f(\lambda) + \sum_{i=1}^l L_i f(\lambda + t_i h). \quad (5.13.4)$$

Равенство содержит $N+m+1$ произвольных параметров α, \dots, λ и A_i, \dots, L_i ($i=0, 1, \dots$). Выбором их можно надеяться сделать правило точным для всех алгебраических многочленов степени $N+m$. Ниже будет показано, что за счет выбора параметров такая степень точности действительно может быть достигнута. Сейчас же мы проверим, что она является наивысшей возможной.

Введем следующие многочлены:

$$\begin{aligned} \omega(x) &= (x-\alpha)(x-\beta)\dots(x-\lambda), \\ \omega_\alpha(x) &= \prod_{i=1}^a (x-\alpha-p_i h), \dots, \omega_\lambda(x) = \prod_{i=1}^l (x-\lambda-t_i h), \\ \Omega(x) &= \omega_\alpha(x)\dots\omega_\lambda(x). \end{aligned}$$

Теорема 1. Ни при каких $\alpha, \beta, \dots, \lambda, A_i, \dots, L_i$ ($i=0, 1, \dots$) правило (5.13.4) не может быть точным для всех многочленов степени $N+m+1$.

Доказательство. Для доказательства достаточно проверить, что при $f=\Omega(x)\omega^2(x)$ равенство (5.13.4) не может выполняться точно. $\Omega\omega^2$ есть многочлен степени $N+m+1$. Узлы квадратурной суммы, стоящей справа в (5.13.4), для него являются

корнями, и сумма равна нулю. С другой стороны, так как многочлен $f = \Omega \omega^2$ сохраняет знак на отрезке $[x_n, x_{n+1}]$ и отличен от тождественного нуля, интеграл

$$\int_{x_n}^{x_{n+1}} \Omega \omega^2 dx$$

не равен нулю и правило (5.13.4) не может быть точным.

Из доказанной теоремы следует, что алгебраическая степень точности (5.13.4) ниже, чем $N+m+1$, и самое большее может быть равна $N+m$.

Теорема 2. Чтобы правил (5.13.4) было точным для всех многочленов степени $N+m$, необходимо и достаточно выполнение условий:

- 1) правило должно быть интерполяционным;
- 2) для всякого многочлена $Q(x)$, степень которого меньше m , должно выполняться равенство

$$\int_{x_n}^{x_{n+1}} \Omega(x) \omega(x) Q(x) dx = 0. \quad (5.13.5)$$

Доказательство. Необходимость первого условия следует из того, что если правило (5.13.4) является точным для многочленов степени $N+m$, то оно точно и для многочленов степени N , а тогда оно должно быть интерполяционным.

Чтобы доказать необходимость второго условия, положим $f = \Omega(x) \omega(x) Q(x)$. Это есть многочлен, степень которого не больше $N+m$, и для него равенство (5.13.4) должно выполняться точно. Но во всех узлах правила f обращается в нуль, и поэтому должно быть верным (5.13.5).

Для проверки достаточности условий положим, что f есть произвольный многочлен степени $N+m$. Если f разделить на $\Omega \omega$, можно f представить в виде

$$f = \Omega \omega Q + r,$$

где r и Q есть многочлены степеней меньше N и m соответственно. Так как $\Omega \omega$ обращается в нуль во всех узлах правила, f и r принимают там одинаковые значения.

Ввиду того что выполняется (5.13.5), а степень r не больше N и правило — интерполяционное, верны следующие равенства, устанавливающие точность (5.13.4) для f :

$$\begin{aligned} \int_{x_n}^{x_{n+1}} f(t) dt &= \int_{x_n}^{x_{n+1}} \Omega \omega Q dt + \int_{x_n}^{x_{n+1}} r dt = \int_{x_n}^{x_{n+1}} r dt = A_0 r(\alpha) + \\ &+ \sum_{i=1}^a A_i r(\alpha + p_i h) + \dots = A_0 f(\alpha) + \sum_{i=1}^a A_i f(\alpha + p_i h) + \dots \end{aligned}$$

Доказанная теорема приводит вопрос о существовании правила, имеющего наивысшую степень точности $N+m$, к вопросу о возможности такого выбора чисел $\alpha, \beta, \dots, \lambda$, чтобы соответствующие им многочлены $\omega(x)$ и $\Omega(x)$ обладали свойством ортогональности (5.13.5).

Теорема 3. Как бы ни были выбраны числа $p_1, \dots, p_i, \dots, p_\lambda$, существуют основные узлы α, \dots, λ , лежащие внутри $[x_n, x_{n+1}]$, для которых выполняется условие (5.13.5).

Доказательство. Возьмем произвольные числа $\alpha, \beta, \dots, \lambda$, удовлетворяющие неравенствам

$$x_n \leq \alpha \leq \beta \leq \dots \leq \lambda \leq x_{n+1}. \quad (5.13.6)$$

По ним составим многочлены $\omega(x)$ и $\Omega(x)$. Многочлен $\Omega(x)$ сохраняет знак на $[x_n, x_{n+1}]$, так как нули его лежат вне или на концах этого отрезка. Примем $\Omega(x)$ за весовую функцию и рассмотрим соответствующую такому весу систему многочленов $P_k(x)$, ортогональную на $[x_n, x_{n+1}]$. Среди них выберем многочлен $P_m(x)$ степени m . Можно считать, что его старший коэффициент равен единице:

$$P_m(x) = x^m + p_1 x^{m-1} + \dots + p_m.$$

Для всякого многочлена $Q(x)$, степень которого меньше m , выполняется равенство

$$\int_{x_n}^{x_{n+1}} \Omega(x) P_m(x) Q(x) dx = 0. \quad (5.13.7)$$

Корни $P_m(x)$ все действительны, различны и лежат внутри отрезка $[x_n, x_{n+1}]$ (см. § 5.5, теорема 3). Обозначим их ξ_k ($k=1, 2, \dots, m$), и пусть $x_1 < \xi_1 < \xi_2 < \dots < \xi_m < x_{n+1}$. Если бы оказалось, что $\xi_1 = \alpha, \xi_2 = \beta, \dots, \xi_m = \lambda$, то многочлен P_m совпал бы с $\omega(x)$ и равенство (5.13.7) было бы равносильно условию (5.13.5).

Покажем, что корни ξ_k непрерывно зависят от $\alpha, \beta, \dots, \lambda$. Свойство ортогональности (5.13.7) эквивалентно равенствам

$$\int_{x_n}^{x_{n+1}} \Omega(x) P_m(x) x^i dx = 0 \quad (i=0, 1, \dots, m-1).$$

Если сюда внести вместо P_m его разложение по степеням x , для коэффициентов p_k получится система уравнений

$$c_{m+i} + c_{m+i-1} p_1 + \dots + c_i p_m = 0 \quad (i=0, 1, \dots, m-1), \quad c_i = \int_{x_n}^{x_{n+1}} \Omega(x) x^i dx.$$

Весовая функция $\Omega(x)$, а следовательно, и числа c_k есть многочлены от $\alpha, \beta, \dots, \lambda$. Определитель системы

$$D = \begin{vmatrix} c_0 & c_1 & \dots & c_{m-1} \\ c_1 & c_2 & \dots & c_m \\ \dots & \dots & \dots & \dots \\ c_{m-1} & c_m & \dots & c_{2m-2} \end{vmatrix}$$

является вместе с тем определителем положительной квадратичной формы переменных z_i ($i=1, 2, \dots, m$)

$$\Phi_m(z_1, z_2, \dots, z_m) = \int_{x_n}^{x_{n+1}} \Omega(x) \left(\sum_{i=1}^m x^{i-1} z_i \right)^2 dx$$

и отличен от нуля при всяких $\alpha, \beta, \dots, \lambda$ из области (5.13.6). Поэтому коэффициенты p_k ($k=1, \dots, m$) есть рациональные функции от $\alpha, \beta, \dots, \lambda$, непрерывные в указанной области.

Корни ξ_k ($k=1, \dots, m$) многочлена $P_m(x)$ непрерывно зависят от p_k и являются поэтому непрерывными функциями α, \dots, λ в той же области:

$$\left. \begin{aligned} \xi_1 &= \varphi_1(\alpha, \beta, \dots, \lambda), \\ \xi_2 &= \varphi_2(\alpha, \beta, \dots, \lambda), \\ &\vdots \\ \xi_m &= \varphi_m(\alpha, \beta, \dots, \lambda). \end{aligned} \right\} \quad (5.13.8)$$

Эти равенства имеют указываемый ниже геометрический смысл. Неравенства (5.13.6) выделяют в m -мерном числовом пространстве замкнутую m -мерную пирамиду. Так как для чисел ξ_k верны соотношения $x_n < \xi_1 < \xi_2 < \dots < \xi_m < x_{n+1}$, равенства (5.13.8) дают однозначное и непрерывное преобразование указанной пирамиды в себя. По теореме Брауэра,* при преобразовании такого рода существует неподвижная точка и, следовательно, существуют значения $\alpha, \beta, \dots, \lambda$, удовлетворяющие неравенствам $x_n < \alpha < \beta < \dots < \lambda < x_{n+1}$, для которых $\xi_1 = \alpha, \xi_2 = \beta, \dots, \xi_m = \lambda$, и будет поэтому выполняться (5.13.5). Из существования же таких $\alpha, \beta, \dots, \lambda$, по теореме 2, вытекает существование правила (5.13.4), имеющего наивысшую степень точности $N+m$.

Вопросы о единственности или числе таких правил не выяснены.

Частные правила такого рода, узлы и коэффициенты для них, так же как соответствующие им остатки, приведены в книгах [3, 4].

5.13.3. О правилах, использующих в вычислениях несколько предшествующих значений интеграла

До сих пор мы рассматривали правила вычисления, в которых для нахождения следующего значения y_{n+1} неопределенного интеграла используется только одно предшествующее его значение. В достаточно общей форме правила такого вида могут быть записаны в виде

$$y_{n+1} \approx y_{n-p} + h \sum_{j=1}^m B_{nj} f(\xi_{nj}) + r_n. \quad (5.13.9)$$

Соответствующее алгебраическое уравнение (5.12.15) имеет форму $\lambda^{p+1}=1$. Корни его по модулю равны единице и являются однократными. Поэтому правило (5.13.9) устойчиво относительно роста погрешности и допускает сходящийся вычислительный процесс всякий раз, когда погрешность r_n правила будет малой величиной выше первого порядка ** сравнительно с h .

Более сложной является теория вычислительных правил, в которых для нахождения y_{n+1} используются несколько предшествующих значений интеграла. Одним из затруднений, возникающих при построении таких правил, является то обстоятельство, что стремление к достижению высокой или даже наивысшей степени точности здесь часто оказывается несовместимым с требованием устойчивости вычислительной схемы.

Поясним эту мысль простым примером. Возвратимся к задаче вычисления неопределенного интеграла от функции, заданной таблицей значений в равноотстоящих точках

* При однозначном и непрерывном отображении m -мерной пирамиды в себя существует по меньшей мере одна неподвижная точка. (См., например, Л. В. Канторович, Г. П. Акилов. Функциональный анализ в нормированных пространствах. М., 1959, стр. 571.)

**) $\frac{r_n}{h} \rightarrow 0$ при $h \rightarrow 0$ равномерно относительно n .

$x_n = x_0 + nh$, и предположим, что вычисления доведены до значения y_n . Допустим теперь, что для вычисления y_{n+1} мы хотим воспользоваться $k+1$ предшествующими значениями функции $y(x)$ и производной от нее: $y_n, y_{n-1}, \dots, y_{n-k}, y'_n = f_n, y'_{n-1} = f_{n-1}, \dots, y'_{n-k} = f_{n-k}$ и построить правило вычислений вида

$$y_{n+1} \approx \sum_{i=0}^k [A_i y_{n-i} + h B_i f_{n-i}]. \quad (5.13.10)$$

Сделаем одно замечание. В задаче неопределенного интегрирования правило такого вида хотя и является допустимым, но не будет иметь наилучшую форму при вычислениях в середине таблицы, так как не использует значений f в точках, следующих за x_n , где f есть известная функция. Правило имеет хорошую форму для вычислений вблизи конца расчетной таблицы, когда мы не знаем значений f дальше X .

Для наших целей такой недостаток формы правила не имеет принципиального значения, и нам достаточно знать, что правило является допустимым и может быть в некоторых случаях полезным. Но, вероятно, следует отметить, что в других, более сложных задачах правила вида (5.13.10) являются естественными и ими широко пользуются в вычислениях. Рассмотрим проблему численного решения дифференциального уравнения первого порядка с начальным условием

$$y' = f(x, y), \quad y(x_0) = x_0, \quad x_0 \leq x \leq X. \quad (5.13.11)$$

Она равносильна решению интегрального уравнения

$$y(x) = y_0 + \int_{x_0}^x f(t, y(t)) dt,$$

и это делает очевидной ее связь с задачей неопределенного интегрирования. Но между этими задачами есть существенное различие. При вычислении неопределенного интеграла (5.13.1) функция f считается известной на всем отрезке $[x_0, X]$ и можно пользоваться любыми значениями f . При численном же нахождении значений решения дифференциального уравнения на сетке точек $x_k = x_0 + kh$ ($k=0, 1, \dots, N$) мы можем считать функцию f известной только до точки x_n и не имеем права без введения дополнительных вычислительных средств пользоваться значениями f в точках, следующих за x_n .

В этой задаче (5.13.11) правило вида (5.13.10), не использующее значений f в точках за x_n , является естественным и его испытание на устойчивость имеет значение.

Правило содержит $2k+2$ произвольных параметров A_i, B_i ($i=0, 1, \dots, k$), и их можно выбрать так, чтобы оно было точным всякий раз, когда $y(x)$ есть многочлен степени не выше $2k+1$. Для этого достаточно сделать правило интерполяционным, при этом интерполирование y_{n+1} здесь выполняется по $2k+2$ значениям функции и первой производной: y_{n-i}, f_{n-i} ($i=0, 1, \dots, k$). Это есть интерполирование с $k+1$ двойными узлами. Нужно нам выражение для y_{n+1} может быть получено при помощи формулы (4.7.10), если в ней положить $x = x_{n+1}$ и внести надлежащие изменения

$$y_{n+1} \approx \sum_{i=0}^k \frac{\omega^2(x_{n+1})}{(x_{n+1} - x_{n-i})^2 \omega'^2(x_{n-i})} \left\{ \left[1 - \frac{\omega''(x_{n-i})}{\omega'(x_{n-i})} (x_{n+1} - x_{n-i}) \right] y_{n-i} + \right. \\ \left. + (x_{n+1} - x_{n-i}) f_{n-i} \right\}, \quad \omega(x) = (x - x_n)(x - x_{n-1}) \dots (x - x_{n-k}). \quad (5.13.12)$$

Как выяснялось в § 5.12, устойчивость или неустойчивость правила (5.13.10) относительно роста погрешности зависит от корней алгебраического уравнения

$$P(\lambda) = \lambda^{k+1} - \sum_{i=0}^k A_i \lambda^{k-i} = 0.$$

Если среди корней есть большие единицы по модулю, то правило будет неустойчивым, при этом неустойчивость будет тем более сильной, чем больше будет максимальный модуль корня. Составить такое уравнение для правила (5.13.12) легко, но запись его в общем виде недостаточно показательна, так как по ней просто судить о наличии корней, имеющих модуль, больший единицы, но затруднительно получить сведения о численных значениях модулей корней.

Рассмотрим коэффициент A_k при нулевой степени λ в многочлене $P(\lambda)$. Абсолютное значение его равно произведению модулей всех корней.

Без труда вычисляются следующие величины:

$$\begin{aligned}\omega(x_{n+1}) &= (x_{n+1} - x_n) \cdot \dots \cdot (x_{n+1} - x_{n-k}) = (k+1)! h^{k+1}, \\ \omega'(x_{n-k}) &= (x_{n-k} - x_n) \cdot \dots \cdot (x_{n-k} - x_{n-k+1}) = (-1)^k k! h^k, \\ \frac{\omega''(x_{n-k})}{\omega'(x_{n-k})} &= 2 \left\{ \frac{1}{x_{n-k} - x_n} + \frac{1}{x_{n-k} - x_{n-1}} + \dots + \frac{1}{x_{n-k} - x_{n-k+1}} \right\} = \\ &= -\frac{2}{h} \sum_{v=1}^k \frac{1}{v}.\end{aligned}$$

Это дает

$$\begin{aligned}A_k &= \frac{\omega^2(x_{n+1})}{(x_{n+1} - x_{n-k})^2 [\omega'(x_{n-k})]^2} \left[1 - \frac{\omega''(x_{n-k})}{\omega'(x_{n-k})} (x_{n+1} - x_{n-k}) \right] = \\ &= 1 + 2(k+1) \sum_{v=1}^k \frac{1}{v}.\end{aligned}$$

При всяких значениях k ($k \geq 1$) A_k больше единицы и среди корней уравнения $P(\lambda) = 0$ всегда есть по модулю большие единицы. Поэтому правило (5.13.12) при $k \geq 1$ всегда неустойчиво в смысле роста погрешностей. Чтобы судить о том, насколько сильной является эта неустойчивость, достаточно рассмотреть частные случаи.

Пусть $k=1$. Интерполирование y_{n+1} будет выполняться по двум двойным узлам x_n, x_{n-1} . Правило (5.13.12) имеет вид $y_{n+1} = -4y_n + 5y_{n-1} + h(4f_n + 2f_{n-1})$. Многочлен $P(\lambda) = \lambda^2 + 4\lambda - 5$. Корни его есть $\lambda_1 = 1, \lambda_2 = -5$. Этим правилом мы пользовались при вычислении интеграла $y(x) = \int_0^x e^t dt$ [см. (5.12.2)] в § 5.12 и видели, насколько сильный

рост погрешности вызывает корень $\lambda_2 = -5$.

Пусть $k=2$. При интерполировании y_{n+1} используются три двойных узла x_n, x_{n-1}, x_{n-2} . Равенство (5.13.12) будет

$$y_{n+1} = -18y_n + 9y_{n-1} + 10y_{n-2} + h(9f_n + 18f_{n-1} + 3f_{n-2}).$$

Корни соответствующего ему многочлена $P(\lambda) = \lambda^3 + 18\lambda^2 - 19\lambda - 10$ есть

$$\lambda_1 = 1, \quad \lambda_2 = \frac{1}{2} (-19 + \sqrt{321}) \approx -0,542, \quad \lambda_3 = \frac{1}{2} (-19 - \sqrt{321}) \approx -18,458.$$

Корень λ_3 позволяет утверждать, что рассматриваемое правило имеет очень быстро возрастающую погрешность и, по-видимому, является малоприменимым для вычислений даже на «небольшое» число шагов.

Сейчас мы ознакомимся с одним из способов составления вычислительных правил рассматриваемого вида, когда требование достижения наивысшей степени точности может быть совмещено с устойчивостью правила. Это достигается, как видно будет ниже, за счет введения «дополнительных» узлов.

Предварительно нам потребуются некоторые сведения из теории интерполирования.

Пусть $y(x)$ есть произвольная, непрерывно дифференцируемая на конечном отрезке $[a, b]$ функция. Возьмем на $[a, b]$ $r+s+l$ разных точек:

$$\begin{aligned} \xi_1, \xi_2, \dots, \xi_r, \\ \xi_{r+1}, \xi_{r+2}, \dots, \xi_{r+s}, \\ \xi_{r+s+1}, \xi_{r+s+2}, \dots, \xi_{r+s+l} \end{aligned}$$

и предположим, что в точках ξ_i ($1 \leq i \leq r$) заданы значения функции $y(\xi_i)$ ($1 \leq i \leq r$); в точках ξ_i ($r < i \leq r+s$) — значения функции и производной $y(\xi_i)$, $y'(\xi_i)$ ($r < i \leq r+s$) и, наконец, в точках ξ_i ($r+s < i \leq r+s+l$) — значения производной $y'(\xi_i)$ ($r+s < i \leq r+s+l$). Для сокращения записи узлы этих трех видов будем называть соответственно простыми, двойными и дополнительными.

Будем считать также, что нам дана на отрезке $[a, b]$ точка x , отличная от двойных и простых узлов. Мы не исключаем заранее возможности ее совпадения с одним из дополнительных узлов, хотя, как видно будет через несколько строк, эта возможность в задаче, которой мы будем заниматься, не может осуществляться.

Поставим теперь себе целью вычислить $y(x)$ по известным значениям $y(\xi_i)$ и $y'(\xi_i)$. Для этого возьмем произвольно $r+s$ чисел α_j ($j=1, 2, \dots, r+s$) и $s+l$ чисел β_j ($j=r+1, \dots, r+s+l$) и составим приближенное равенство

$$y(x) \approx \sum_{j=1}^{r+s} \alpha_j y(\xi_j) + \sum_{j=r+1}^{r+s+l} \beta_j y'(\xi_j). \quad (5.13.13)$$

Алгебраическая степень точности ^{*)} равенства зависит как от коэффициентов α_j , β_j , так и от точек x , ξ_j ($j=1, \dots, r+s+l$). Выясним, какой наивысшей возможной степени точности может достигнуть равенство и каким условиям должны для этого удовлетворять параметры α_j , β_j , x , ξ_j .

Теорема 4. При всяких α_j , β_j и любом положении точек x , ξ_j степень точности правила (5.13.13) всегда меньше $r+2s+2l$.

Доказательство. Достаточно показать, что всегда существует многочлен, имеющий степень не больше $r+2s+2l$, для которого равенство (5.13.13) не может выполняться точно.

Пусть x не совпадает ни с одним из дополнительных узлов.

Рассмотрим многочлен степени $r+2s+2l$

$$y(z) = (z - \xi_1) \dots (z - \xi_r) (z - \xi_{r+1})^2 \dots (z - \xi_{r+s+1})^2 = A(z). \quad (5.13.14)$$

Сразу же видно, что $A(\xi_j) = 0$ ($j=1, \dots, r+s+l$) и $A'(\xi_j) = 0$ ($j > r$). Для $y(z) = A(z)$ правая часть (5.13.13) равна нулю, тогда как $y(x) = A(x) \neq 0$ и равенство (5.13.13) не может быть точным.

^{*)} Напомним, что число m называют степенью точности равенства, если равенство выполняется точно для всякого многочлена $y(z)$ степени m и не является точным для $y(z) = z^{m+1}$.

Допустим теперь, что x совпадает с одним из дополнительных узлов, например, пусть $x = \xi_{r+s+l}$, и введем многочлен $B(z) = \frac{A(z)}{(z-x)^2}$. При $B'(\xi_{r+s+l}) = B'(x) \neq 0$ положим

$$y(z) = B(z) \left[z - x - \frac{B(x)}{B'(x)} \right].$$

$y(z)$ есть многочлен степени $r+2s+2l+1$. Ввиду $y'(\xi_{r+s+l}) = y'(x) = 0$ правая часть (5.13.13) равна нулю, левая же часть

$$y(x) = -\frac{B^2(x)}{B'(x)}$$

отлична от нуля и равенство (5.13.13) не может быть точным. Поэтому в рассматриваемом случае степень точности (5.13.13) ниже $r+2s+2l-1$.

Когда $B'(x) = B'(\xi_{r+s+l}) = 0$, достаточно положить $y(z) = B(z)$, чтобы убедиться в том, что степень точности (5.13.13) в этом случае меньше $r+2s+2l-2$.

Теорема 4 показывает, что степень точности m равенства (5.13.13) всегда меньше $r+2s+2l$ и самое большее может быть равна $r+2s+2l-1$. Из доказательства теоремы видно, что если $m = r+2s+2l-1$, то точка x должна быть отличной от всех дополнительных узлов ξ_j ($r+s < j \leq r+s+l$). Это условие всюду ниже будем считать выполненным.

Докажем теорему, дающую условия, при которых равенство (5.13.13) имеет наивысшую степень точности $m = r+2s+2l-1$.

Теорема 5. Для того чтобы равенство (5.13.13) было точным для всяких многочленов степени $r+2s+2l-1$, необходимо и достаточно выполнение условий:

1) x и ξ_j должны удовлетворять системе l уравнений

$$\sum'_{k=r+s+1}^{r+s+l} \frac{2}{\xi_j - \xi_k} + \sum_{k=1}^r \frac{2}{\xi_j - \xi_k} + \sum_{k=r+1}^{r+s} \frac{2}{\xi_j - \xi_k} + \frac{1}{\xi_j - x} = 0 \quad (5.13.15)$$

$$(j = r+s+1, \dots, r+s+l),$$

(знак штрих, стоящий у Σ , показывает, что должно быть пропущено значение $k=j$);

2) коэффициенты α_j и β_j должны иметь значения:

$$\begin{aligned} \alpha_j &= \frac{A(x)}{(x - \xi_j)A'(\xi_j)} \quad (j = 1, 2, \dots, r), \\ \alpha_j &= \frac{A_j(x)}{A_j(\xi_j)} \left[1 - (x - \xi_j) \frac{A_j'(\xi_j)}{A_j(\xi_j)} \right] \quad (j = r+1, \dots, r+s), \\ \beta_j &= -\frac{A_j(x)}{A_j(\xi_j)} (x - \xi_j) \quad (j = r+1, \dots, r+s+l), \end{aligned} \quad (5.13.16)$$

где

$$A_j(z) = A(z)(z - \xi_j)^{-2}.$$

Перед доказательством сделаем замечание о теореме. Она приводит задачу о построении равенства (5.13.13), верного для многочленов степени $r+2s+2l-1$, к вопросу о выполнении уравнений (5.13.15), так как по x и ξ_j коэффициенты α_j и β_j всегда могут быть построены единственным образом при помощи формул (5.13.16).

Построить простое доказательство теоремы можно, если воспользоваться правилами интерполирования с кратными узлами (§ 4.7).

Доказательство. Пусть $y(z)$ есть произвольная дифференцируемая функция на $[a, b]$. Интерполируем ее по следующим значениям $y(z)$ и $y'(z)$:

$$y(\xi_j) \quad (j=1, 2, \dots, r+s+l), \quad y'(\xi_j) \quad (j=r+1, \dots, r+s+l).$$

Это есть интерполирование с r простыми узлами ξ_j ($j=1, \dots, r$) и $s+l$ двойными узлами ξ_j ($r < j \leq r+s+l$). Интерполирующий многочлен $P(y; z)$, имеющий степень не больше $r+2s+2l-1$, может быть построен по правилу (4.7.8) и имеет следующий вид:

$$\begin{aligned} P(z) = P(y; z) &= \sum_{k=1}^r \frac{A(z)}{(z-\xi_k)A'(\xi_k)} y(\xi_k) + \\ &+ \sum_{k=r+1}^{r+s+l} \frac{A_k(z)}{A_k(\xi_k)} \left[1 - (z-\xi_k) \frac{A_k'(\xi_k)}{A_k(\xi_k)} \right] y(\xi_k) + \sum_{k=r+1}^{r+s+l} \frac{A_k(z)}{A_k(\xi_k)} (z-\xi_k) y'(\xi_k) = \\ &= \sum_{k=1}^{r+s+l} l_k(z) y(\xi_k) + \sum_{k=r+1}^{r+s+l} d_k(z) y'(\xi_k). \end{aligned} \quad (5.13.17)$$

Остаток интерполирования $R(y; z) = y(z) - P(y; z)$ имеет значение *)

$$R(y; z) = \frac{A(z)}{(r+2s+2l)!} y^{(r+2s+2l)}(\xi).$$

Наиболее просто убедиться в справедливости (5.13.17) можно при помощи простой проверки равенств

$$P(\xi_j) = y(\xi_j) \quad (j=1, \dots, r+s+l) \quad \text{и} \quad P'(\xi_j) = y'(\xi_j) \quad (j=r+1, \dots, r+s+l).$$

Коэффициенты $l_k(z)$ и $d_k(z)$ являются многочленами влияния соответствующих значений $y(\xi_k)$ или $y'(\xi_k)$. Каждый из них имеет степень $r+2s+2l-1$.

При $k > r$ многочлен

$$l_k(z) = \frac{A_k(z)}{A_k(\xi_k)} \left[1 - (z-\xi_k) \frac{A_k'(\xi_k)}{A_k(\xi_k)} \right]$$

обладает следующими свойствами:

$$l_k(\xi_j) = \begin{cases} 0, & j \neq k, \\ 1, & j = k, \end{cases} \quad (j=1, \dots, r+s+l), \quad l_k'(\xi_j) = 0 \quad (j=r+1, \dots, r+s+l). \quad (5.13.18)$$

Аналогичные свойства легко указать и для других многочленов влияния $l_k(z)$ ($k=1, \dots, r$) и $d_k(z)$ ($k > r$).

Проверим необходимость условий теоремы. Пусть равенство (5.13.13) точно для многочленов степени $r+2s+2l-1$. В частности, оно должно быть точным для многочлена $y(z) = l_j(z)$ при $j > r+s$. Но правая часть (5.13.13) для такого многочлена влияния равна нулю, и должно, следовательно, выполняться равенство

*) Приведенное представление остатка R предполагает существование у $y(x)$ непрерывной производной порядка $r+2s+2l$ на $[a, b]$.

$$l_j(x) = \frac{A_j(x)}{A_j(\xi_j)} \left[1 - (x - \xi_j) \frac{A_j'(\xi_j)}{A_j(\xi_j)} \right] = 0 \quad (j = r+s+1, \dots, r+s+l). \quad (5.13.19)$$

Если сократить на не равный нулю множитель

$$-\frac{A_j(x)}{A_j(\xi_j)} (x - \xi_j)$$

и вычислить логарифмическую производную $\frac{A_j'(\xi_j)}{A_j(\xi_j)}$, сразу же видно будет, что эти равенства совпадут с (5.13.15). Этим доказана необходимость первого условия теоремы.

Чтобы доказать необходимость выполнения первого из равенств, входящих во второе условие теоремы, будем считать j лежащим в границах от 1 до r и рассмотрим многочлен влияния

$$l_j(z) = \frac{A(z)}{(z - \xi_j) A'(\xi_j)}.$$

Он имеет степень $r+2s+2l-1$ и для него равенство (5.13.13) выполняется точно. Но

$$l_j(\xi_k) = \begin{cases} 1, & k=j, \\ 0, & k \neq j, \end{cases} \quad (k=1, 2, \dots, r+s+l), \quad l_j'(\xi_k) = 0 \quad (k=r+1, \dots, r+s+l)$$

и (5.13.13) примет вид

$$l_j(x) = \frac{A(x)}{(x - \xi_j) A'(\xi_j)} = \alpha_j \cdot 1.$$

Аналогично доказываются и другие равенства (5.13.16).

Достаточность условий проверяется столь же просто. Если $y(z)$ есть произвольный многочлен степени $r+2s+2l-1$, для него интерполяционный многочлен (5.13.17) будет совпадать с $y(z)$ при всяком z . В частности, в точке x будет верно равенство

$$\begin{aligned} y(x) = & \sum_{k=1}^r \frac{A(x)}{(x - \xi_k) A'(\xi_k)} y(\xi_k) + \sum_{k=r+1}^{r+s+l} \frac{A_k(x)}{A_k(\xi_k)} \left[1 - (x - \xi_k) \frac{A_k'(\xi_k)}{A_k(\xi_k)} \right] y(\xi_k) + \\ & + \sum_{k=r+1}^{r+s+l} \frac{A_k(x)}{A_k(\xi_k)} (x - \xi_k) y'(\xi_k). \end{aligned} \quad (5.13.20)$$

Условия (5.13.15) равносильны равенствам (5.13.19) и в средней сумме $\sum_{k=r+1}^{r+s+l}$ верхний индекс суммирования $r+s+l$ может быть заменен на $r+s$. А тогда, ввиду значений (5.13.16) коэффициентов α_j и β_j , последнее равенство становится равносильным (5.13.13) и, так как это верно для всякого многочлена $y(z)$ степени $r+2s+2l-1$, достаточность можно считать установленной.

Теорема 5 позволяет сказать, что мы можем достигнуть в равенстве (5.13.13) наивысшей степени точности $r+2s+2l-1$, если узлы ξ_j выберем так, чтобы выполнялись равенства (5.13.15). После этого достаточно коэффициентам α_j , β_j придать значения (5.13.16).

Мы должны теперь выяснить, можно ли выполнить условия (5.13.15), и если можно, то какой произвол останется после этого в нашем распоряжении. Большую помощь нам окажет то обстоятельство, что условия имеют простое и очень наглядное физическое истолкование. Рассмотрим комплексную плоскость, возьмем в ней две точки z_1, z_2 и поместим в них электрические заряды масс e_1 и e_2 . Предположим, что они действуют друг на друга с силой, численное значение которой пропорционально массам e_1, e_2 и обратно пропорционально первой степени расстояния, направленной по прямой, соединяющей z_1 и z_2 . С взаимодействием такого рода приходится иметь дело в теории плоского электростатического поля.

Коэффициент пропорциональности можно принять равным единице. Сила отталкивания, с которой заряд в точке z_1 действует на заряд в z_2 , равна $\frac{e_1 e_2}{z_2 - z_1}$.

Возьмем теперь на плоскости $r+s+1$ точек $x, \xi_1, \xi_2, \dots, \xi_{r+s}$ и закрепим их. В точки x, ξ_1, \dots, ξ_r поместим заряды массы 1, а в точки $\xi_{r+1}, \dots, \xi_{r+s}$ — заряды массы 2. Кроме того, возьмем l свободных зарядов массы 2 и комплексные координаты их обозначим $\xi_{r+s+1}, \dots, \xi_{r+s+l}$. В положении равновесия равнодействующие всех сил, приложенных к каждому свободному заряду, должны быть равны нулю:

$$\sum_{k=r+s+1}^{r+s+l} \frac{2 \cdot 2}{\xi_j - \xi_k} + \sum_{k=1}^r \frac{2 \cdot 1}{\xi_j - \xi_k} + \sum_{k=r+1}^{r+s} \frac{2 \cdot 2}{\xi_j - \xi_k} + \frac{2 \cdot 1}{\xi_j - x} = 0$$

$$(j = r+s+1, \dots, r+s+l).$$

Эти уравнения равносильны (5.13.15) и отличаются от них только множителем 2 и сопряжением.

Указанная сейчас аналогия между условиями (5.13.15) и положениями равновесия системы электрических зарядов позволяет исследовать уравнения (5.13.15) наглядным путем и сделать очевидными указываемые ниже результаты.

1. Если x, ξ_j ($j=1, \dots, r+s$) есть любые комплексные числа и ξ_j ($j=r+s+1, \dots, r+s+l$) удовлетворяют системе (5.13.15), то точки с комплексными координатами $\xi_{r+s+1}, \dots, \xi_{r+s+l}$ лежат в наименьшем выпуклом многоугольнике, содержащем x и ξ_j ($j \leq r+s$). В частном случае, когда $x, \xi_1, \dots, \xi_{r+s}$ лежат на действительной оси, то ξ_j ($j > r+s$) лежат внутри наименьшего отрезка, содержащего $x, \xi_1, \dots, \xi_{r+s}$.

2. Пусть $x, \xi_1, \dots, \xi_{r+s}$ действительны и различны. На числовой оси между точками с такими координатами будет $r+s$ промежутков. Предположим, что указан закон распределения свободных зарядов по таким промежуткам, т. е. указано, сколько свободных зарядов должно лежать в каждом из промежутков. Отметим, что число способов рас-

пределения можно подсчитать и оно равно $\frac{(r+s+l-1)!}{l!(r+s-1)!}$. Существует решение системы

(5.13.15), имеющее заданный закон распределения свободных зарядов по промежуткам, и если не различать решения, получающиеся друг из друга перестановкой свободных зарядов, то такое решение будет единственным.

Изложенные простые соображения позволяют высказать приводимую ниже теорему о равенстве (5.13.13).

Теорема 6. При всяких действительных и различных $x, \xi_1, \dots, \xi_{r+s}$ дополнительные узлы $\xi_{r+s+1}, \dots, \xi_{r+s+l}$ всегда можно выбрать $\frac{(r+s+l-1)!}{l!(r+s-1)!}$ способами так, чтобы равенство (5.13.13) имело наивысшую алгебраическую степень точности $r+2s+2l-1$. При этом, если заранее указать, какое число дополнительных узлов должно лежать в каждом из промежутков между $x, \xi_1, \dots, \xi_{r+s}$, то среди возможных способов достижения наивысшей степени точности существует один и только один способ, при котором дополнительные узлы имеют указанное для них распределение по промежуткам.

Возвратимся к задаче вычисления неопределенного интеграла (5.13.1). Предположим, что последнее найденное значение есть y_n , и постараемся выяснить рациональную форму применения правила (5.13.13) к нахождению y_{n+1} . Роль точки интерполирования x здесь играет табличный узел x_{n+1} . Естественно требовать, чтобы правило имело наивысшую степень точности. Напомним, что тогда в нашем распоряжении находится выбор простых и двойных узлов ξ_j ($j=1, \dots, r+s$). В качестве их могут быть взяты любые табличные узлы x_k , предшествующие x_{n+1} . Предположим, что они как-то избраны и фиксированы. После этого мы имеем право распорядиться распределением вспомогательных узлов между точками $x_{n+1}, \xi_1, \dots, \xi_{r+s}$.

Легко предвидеть, какой закон распределения следует признать наилучшим, если строить правило вычислений для широкого класса функций и не учитывать какие-либо особые их свойства, которые могли бы заставить по-разному оценивать различные участки отрезка $[a, b]$. В этих условиях для достижения лучшей точности выгодно вспомогательные узлы взять возможно ближе к точке, где мы вычисляем функцию $y(x)$, и поместить их в промежуток, примыкающий к x_{n+1} .

Можно привести простые соображения, подтверждающие такое предвидение, связанные с оценкой погрешности правила (5.13.13). При доказательстве теоремы 5 мы обращали внимание на то, что правая часть равенства (5.13.13), если оно имеет наивысшую степень точности, совпадает со значением вспомогательного интерполяционного многочлена $P(z) = P(y; z)$ (5.13.17) при $z=x$. Поэтому погрешность правила равна значению остатка интерполирования $R(f; x)$, аналитическое представление которого приведено двумя строками ниже (5.13.17):

$$R(y; x) = \frac{A(x)}{(r+2s+2)!} y^{(r+2s+2)}(\xi),$$

$$A(x) = (x-\xi_1) \dots (x-\xi_r) (x-\xi_{r+1})^2 \dots (x-\xi_{r+s+1})^2.$$

Для множества функций, определяемого условием

$$|y^{(r+2s+2)}(z)| \leq M, \quad z \in [a, b],$$

погрешность правила имеет следующую точную оценку:

$$|R(y; x)| \leq \frac{|A(x)|}{(r+2s+2)!} M.$$

От выбора узлов ξ_j в оценке зависит только $|A(x)|$, и наилучшим расположением узлов естественно считать то, при котором эта величина имеет наименьшее значение. Если считать простые и двойные узлы ξ_j ($j \leq r+s$) лежащими слева от x_{n+1} и если принять во внимание электростатическую аналогию задачи о расположении узлов, то будет ясно, что множители, входящие в $A(x)$ и отвечающие вспомогательным узлам $x-\xi_{r+s+1}, \dots, \xi_{r+s+1}$, будут иметь наименьшие численные значения, когда $\xi_{r+s+1}, \dots, \xi_{r+s+1}$ будут лежать в промежутке, примыкающем к x_{n+1} .

Произведение $(x_{n+1}-\xi_1) \dots (x_{n+1}-\xi_r) (x_{n+1}-\xi_{r+1})^2 \dots (x_{n+1}-\xi_{r+s})^2$, отвечающее простым и двойным узлам, будет иметь наименьшее значение, очевидно, в том случае, когда эти узлы займут табличные места, ближайшие предшествующие x_{n+1} , при этом двойные узлы следует поместить в точках $x_n, x_{n-1}, \dots, x_{n-s+1}$ и простые — в точках $x_{n-s}, \dots, x_{n-s-r+1}$.

Если воспользоваться электростатической аналогией, то можно сказать, что при таком расположении двойных и простых зарядов они будут оказывать на свободные заряды наиболее сильное давление слева направо и прижимать их к заряду в точке x_{n+1} . Поэтому множители $(x_{n+1}-\xi_j)$ ($j > r+s$), отвечающие свободным зарядам, будут иметь наименьшие возможные значения.

Приведенные наглядные соображения позволяют высказать следующее правило построения расчетной формулы (5.13.13).

Если равенство (5.13.13) имеет наивысшую степень точности и его узлы ξ_j и коэффициенты α_j, β_j удовлетворяют, следовательно, условиям теоремы 5, то для достижения минимального значения оценки погрешности правила нужно:

- 1) дополнительные узлы поместить внутри отрезка $[x_n, x_{n+1}]$;
- 2) двойные узлы взять в точках $x_n, x_{n-1}, \dots, x_{n-s+1}$;
- 3) в качестве простых узлов взять $x_{n-s}, \dots, x_{n-s-r+1}$.

При таком выборе узлов правило (5.13.13) принимает вид

$$y_{n+1} \approx \sum_{j=0}^{s-1} (a_j y_{n-j} + b_j y'_{n-j}) + \sum_{j=0}^{r-1} c_j y_{n-s-j} + \sum_{j=1}^l d_j y'(x_n + ht_j) \quad (5.13.21)$$

$$(0 < t_1 < t_2 < \dots < t_l < 1).$$

Покажем теперь, что правило вида (5.13.21), имеющее наивысшую степень точности, всегда может быть построено так, чтобы оно было устойчивым относительно роста погрешности. Для этого достаточно, как показано в теореме 5 § 5.12, чтобы коэффициенты a_j и c_j при значениях функции были неотрицательны. Чтобы не вводить новых обозначений, возвратимся к старой записи (5.13.13) расчетной формулы, но положим в ней $x = x_{n+1}$ и будем считать, что ее узлы ξ_j в соответствии с указанным несколькими строками выше правилом взяты на отрезке, прилежащем к точке x_{n+1} .

Мы должны выяснить условия положительности коэффициентов α_j .

Возьмем один из простых узлов ξ_j ($1 \leq j \leq r$) и рассмотрим соответствующий ему коэффициент

$$\alpha_j = \frac{A(x_{n+1})}{(x_{n+1} - \xi_j) A'(\xi_j)}, \quad A(x) = (x - \xi_1) \dots (x - \xi_r) (x - \xi_{r+1})^2 \dots (x - \xi_{r+s+1})^2.$$

Так как узлы ξ_j все лежат слева от x_{n+1} , значения $A(x_{n+1})$ и $x_{n+1} - \xi_j$ положительны и знак α_j совпадает со знаком $A'(\xi_j)$.

$$A'(\xi_j) = (\xi_j - \xi_1) \dots (\xi_j - \xi_{j-1}) (\xi_j - \xi_{j+1}) \dots (\xi_j - \xi_r) (\xi_j - \xi_{r+1})^2 \dots (\xi_j - \xi_{r+s+1})^2.$$

Отсюда видно, что значения $A'(\xi_j)$, а следовательно и α_j , отвечающие соседним простым узлам, будут противоположных знаков. Если мы хотим, чтобы все α_j были положительными, необходимо считать либо $r=0$ и не брать ни одного простого узла, либо $r=1$ и взять лишь один простой узел.

Рассмотрим теперь коэффициент α_j ($r < j \leq r+s$), отвечающий двойному узлу:

$$\alpha_j = \frac{A_j(x_{n+1}) (\xi_j - x_{n+1})}{A_j(\xi_j)} \left[\frac{1}{\xi_j - x_{n+1}} + \frac{A_j'(\xi_j)}{A_j(\xi_j)} \right],$$

$$A_j(z) = A(z) (z - \xi_j)^{-2}.$$

Простые узлы лежат левее двойных, и поэтому $A_j(\xi_j) > 0$. Так как $\xi_j < x_{n+1}$, множитель перед прямоугольной скобкой отрицателен и условием неотрицательности α_j будет неравенство

$$\frac{1}{\xi_j - x_{n+1}} + \frac{A_j'(\xi_j)}{A_j(\xi_j)} = \frac{1}{\xi_j - x_{n+1}} + \sum_{k=1}^r \frac{1}{\xi_j - \xi_k} + \sum_{k=r+1}^{r+s+1} \frac{2}{\xi_j - \xi_k} \leq 0. \quad (5.13.22)$$

Ему легко придать физический смысл. Возвратимся к электростатической аналогии условий достижения наивысшей степени точности. В левой части неравенства (5.13.22)

стоит величина, равная половине значения сил отталкиваний, которые испытывает заряд в точке ξ_j со стороны всех других зарядов. Условие же (5.13.22) означает, что эта равнодействующая направлена по числовой оси в отрицательную сторону или равна нулю.

Точка x_{n+1} и все вспомогательные узлы лежат справа от двойного узла и будут толкать заряд в ξ_j налево. Это дает право высказать следующее утверждение.

Какими бы ни были числа r и s , для них существует l_0 такое, что при $l \geq l_0$ все α_j ($j=r+1, \dots, r+s$) будут неотрицательны.

Например, допустим, что берется один двойной узел в точке x_n и один простой узел в точке x_{n-1} . Коэффициент α_1 , отвечающий простому узлу x_{n-1} , как выяснено выше, будет положителен. Условие положительности коэффициента α_2 , отвечающего двойному узлу x_n , запишется в виде неравенства

$$\frac{1}{x_n - x_{n+1}} + \frac{1}{x_n - x_{n-1}} + \sum_{j=3}^{2+l} \frac{2}{x_n - \xi_j} = \sum_{j=3}^{2+l} \frac{2}{x_n - \xi_j} \leq 0.$$

Так как $\xi_j > x_n$, оно выполняется при всяких $l \geq 1$.

Численные значения коэффициентов α_j , β_j и вспомогательных узлов вычислены для некоторых простейших случаев. Небольшие таблицы их можно найти в книгах [3, 4].

Литература

1. Березин И. С., Жидков Н. П. Методы вычислений, т. I. М., 1966.
2. Крылов А. Н. Лекции о приближенных вычислениях. Л., 1933.
3. Крылов В. И. Приближенное вычисление интегралов. М., 1967.
4. Крылов В. И., Шultzгина Л. Т. Справочная книга по численному интегрированию. М., 1966.

Добавление I

НЕКОТОРЫЕ СВЕДЕНИЯ ИЗ ФУНКЦИОНАЛЬНОГО АНАЛИЗА

§ 1. МЕТРИЧЕСКИЕ ПРОСТРАНСТВА. СХОДИМОСТЬ И ПОЛНОТА

Пусть X есть множество элементов x произвольной природы. Оно называется метрическим пространством, если каждому двум элементам $x, y \in X$ поставлено в соответствие число $\rho(x, y)$, называемое расстоянием между x и y и удовлетворяющее условиям:

- 1) $\rho(x, y) \geq 0$, $\rho(x, y) = 0$ в том и только в том случае, когда $x = y$;
- 2) $\rho(x, y) = \rho(y, x)$;
- 3) $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$.

Эти условия являются аксиомами метрики. Первое из них называют иногда аксиомой различения, второе — аксиомой симметрии и третье — аксиомой треугольника.

Расстояние позволяет естественным образом ввести понятие сходимости. Последовательность элементов x_n ($n = 1, 2, \dots$) называется сходящейся к элементу x^* , если $\rho(x_n, x^*) \rightarrow 0$ ($n \rightarrow \infty$). В таком случае пишут либо $x_n \rightarrow x^*$, либо $\lim x_n = x^*$.

Укажем на два простых, но важных следствия из аксиом метрики.

1. Если $x_n \rightarrow x^*$ и $y_n \rightarrow y^*$, то $\rho(x_n, y_n) \rightarrow \rho(x^*, y^*)$, иначе говоря, расстояние $\rho(x, y)$ есть непрерывная функция своих аргументов x и y .

При помощи двукратного применения неравенства треугольника получим соотношение

$$\rho(x', y') \leq \rho(x', x) + \rho(x, y') \leq \rho(x', x) + \rho(x, y) + \rho(y, y').$$

Отсюда

$$\rho(x', y') - \rho(x, y) \leq \rho(x', x) + \rho(y', y).$$

Меняя же местами пары элементов x, y и x', y' , получим неравенство с противоположным знаком в левой части, и поэтому

$$|\rho(x', y') - \rho(x, y)| \leq \rho(x', x) + \rho(y', y).$$

На этом основании можем написать

$$|\rho(x_n, y_n) - \rho(x^*, y^*)| \leq \rho(x_n, x^*) + \rho(y_n, y^*) \rightarrow 0,$$

что доказывает утверждение.

2. Если $x_n \rightarrow x^*$ и $x_n \rightarrow x'$, то $x^* = x'$ и сходящаяся последовательность не может иметь двух разных пределов.

Действительно, $0 = \rho(x_n, x_n) \rightarrow \rho(x^*, x')$. Следовательно, $\rho(x^*, x') = 0$ и $x^* = x'$.

Для всякой сходящейся последовательности элементов x_n выполняется признак Больцано — Коши. В самом деле, если $x_n \rightarrow x$, то для величины $\varepsilon > 0$ существует n_0 такое, что при всяком $n > n_0$ будет $\rho(x_n, x) < \frac{1}{2} \varepsilon$. Пусть $n > n_0$ и $m > n_0$. Тогда

$$\rho(x_n, x_m) \leq \rho(x_n, x) + \rho(x, x_m) \leq \frac{1}{2} \varepsilon + \frac{1}{2} \varepsilon = \varepsilon.$$

Обратное, вообще говоря, не верно, так как существуют метрические пространства, в которых из выполнения признака Больцано — Коши для последовательности не обязательно следует, что эта последовательность будет сходящейся.

Метрическое пространство X называется полным, если в нем всякая последовательность x_n , для которой выполняется признак Больцано — Коши, будет сходящейся. Для полных метрических пространств верна теорема о сходимости: для того, чтобы последовательность элементов x_n была сходящейся, необходимо и достаточно, чтобы для нее выполнялся признак Больцано — Коши.

Пример неполного метрического пространства дает множество R рациональных чисел, в котором за расстояние принято абсолютное значение разности между числами: $\rho(x, y) = |x - y|$. Последовательность рациональных чисел x_n , выполняющая признак Больцано — Коши, может не иметь предела в R , так как ее пределом может быть иррациональное число, которое не принадлежит R .

Приведем примеры метрических пространств, полезных в теории систем численных уравнений.

Рассмотрим n -мерное числовое пространство R_n , элементами которого являются упорядоченные совокупности n действительных (или комплексных) чисел $x = (x_1, x_2, \dots, x_n)$. Метрика в R_n может быть введена многими способами. Мы остановимся на трех наиболее употребительных метриках.

1. Кубическая, или m -метрика. Она определяется равенством

$$\rho_m(x, y) = \max_i |x_i - y_i|. \quad (I.1)$$

Первая и вторая аксиомы метрики, очевидно, выполняются. Выполнение третьей аксиомы проверяется легко:

$$\rho_m(x, z) = \max_i |x_i - z_i| = \max_i |(x_i - y_i) + (y_i - z_i)| \leq$$

$$\begin{aligned} \leq \max_i [|x_i - y_i| + |y_i - z_i|] &\leq \max_i |x_i - y_i| + \max_i |y_i - z_i| = \\ &= \rho_m(x, y) + \rho_m(y, z). \end{aligned}$$

2. Октаэдрическая, или s -метрика. Определяется равенством

$$\rho_s(x, y) = \sum_{i=1}^n |x_i - y_i|.$$

Выполнение первой и второй аксиом метрики является очевидным. Третья аксиома просто проверяется:

$$\begin{aligned} \rho_s(x, z) &= \sum_{i=1}^n |x_i - z_i| = \sum_{i=1}^n |(x_i - y_i) + (y_i - z_i)| \leq \\ &\leq \sum_{i=1}^n |x_i - y_i| + \sum_{i=1}^n |y_i - z_i| = \rho_s(x, y) + \rho_s(y, z). \end{aligned}$$

3. Сферическая, или l -метрика:

$$\rho_l(x, y) = \left\{ \sum_{i=1}^n (x_i - y_i)^2 \right\}^{\frac{1}{2}}.$$

Первая и вторая аксиомы метрики столь же очевидно верны, как и в предыдущих метриках, третья же аксиома есть простое следствие неравенства Коши

$$\sum_{i=1}^n a_i b_i \leq \left\{ \sum_{i=1}^n a_i^2 \right\}^{\frac{1}{2}} \left\{ \sum_{i=1}^n b_i^2 \right\}^{\frac{1}{2}}.$$

Действительно,

$$\begin{aligned} \rho_l^2(x, z) &= \sum_{i=1}^n (x_i - z_i)^2 = \sum_{i=1}^n [(x_i - y_i) + (y_i - z_i)]^2 = \sum_{i=1}^n (x_i - y_i)^2 + \\ &+ \sum_{i=1}^n (y_i - z_i)^2 + 2 \sum_{i=1}^n (x_i - y_i)(y_i - z_i) \leq \rho^2(x, y) + \rho^2(y, z) + \\ &+ 2\rho(x, y)\rho(y, z) = [\rho(x, y) + \rho(y, z)]^2. \end{aligned}$$

Во всех трех приведенных метриках сходимость последовательности элементов $x^{(m)} \rightarrow x^*$ ($m \rightarrow \infty$), очевидно, равносильна n численным «покоординатным» сходимостям $x_i^{(m)} \rightarrow x_i^*$ ($i = 1, 2, \dots, n$). Различными будут лишь «меры скоростей сходимости».

§ 2. ЛИНЕЙНЫЕ НОРМИРОВАННЫЕ ПРОСТРАНСТВА. ЛИНЕЙНЫЕ ОПЕРАТОРЫ

Множество X называется линейным множеством или векторным пространством, если для каждого двух его элементов x и y определена сумма $x+y$, являющаяся элементом того же множества X , и для всякого элемента x и числа λ (действительного или комплексного) определено произведение λx , которое есть элемент X , при этом обе указанные операции подчиняются следующим условиям:

- 1) ассоциативность сложения $(x+y)+z=x+(y+z)$;
- 2) коммутативность сложения $x+y=y+x$;
- 3) существует элемент 0 , называемый нулевым, такой, что при всяком x из X $x+0=x$;
- 4) для каждого x существует элемент $-x$, называемый противоположным, такой, что $x+(-x)=0$;
- 5) ассоциативность умножения $\lambda(\mu x)=(\lambda\mu)x$;
- 6) два распределительных закона: $(\lambda+\mu)x=\lambda x+\mu x$, $\lambda(x+y)=\lambda x+\lambda y$;
- 7) $1 \cdot x=x$;
- 8) $0 \cdot x=0$;
- 9) если $\lambda x=0$ и $x \neq 0$, то $\lambda=0$.

Линейное множество называется линейным нормированным пространством, если для каждого элемента $x \in X$ определена норма $\|x\|$, являющаяся действительным числом и удовлетворяющая условиям:

- 1) $\|x\| \geq 0$ и $\|x\|=0$ тогда и только тогда, когда $x=0$;
- 2) $\|x+y\| \leq \|x\| + \|y\|$;
- 3) $\|\lambda x\| = |\lambda| \cdot \|x\|$.

В линейном нормированном пространстве всегда может быть введена метрика, для чего достаточно за расстояние между элементами x и y принять норму их разности $x-y$:

$$\rho(x, y) = \|x - y\|.$$

Это позволяет определить в X сходимость последовательности элементов: если $\|x_n - x^*\| \rightarrow 0$, то $x_n \rightarrow x^*$, и определить полноту пространства.

Линейное нормированное и полное пространство называется пространством типа Бахана или, коротко, типа В. В таком пространстве всякий «абсолютно» (по норме) сходящийся ряд будет сходиться, т. е. из

сходимости числового ряда $\sum_{k=1}^{\infty} \|x_k\|$ следует существование предела

$$\lim_{n \rightarrow \infty} \sum_{k=1}^n x_k = \sum_{k=1}^{\infty} x_k.$$

Приведем нужные нам примеры пространств типа В.

1. Пространство $C[a, b]$. Пусть $[a, b]$ есть конечный замкнутый отрезок и рассматривается множество функций $x(t)$, непрерывных на $[a, b]$. Сложение элементов и умножение на число есть обычное сложение функций и умножение их на число. За норму функции принимают максимум абсолютного ее значения

$$\|x(t)\| = \max_t |x(t)|.$$

Сходимость последовательности элементов из $[a, b]$ есть равномерная сходимость последовательности функций. Пространство $C[a, b]$ является полным. Из $\|x_{n+m}(t) - x_n(t)\| = \max_t |x_{n+m}(t) - x_n(t)| \leq \varepsilon$ вытекает сходимость $x_n(t)$ при всяком $t \in [a, b]$: $\lim_{n \rightarrow \infty} x_n(t) = x(t)$, и так как предел равномерно сходящейся последовательности непрерывных функций есть функция непрерывная, то $x(t) \in C[a, b]$.

2. Пространство $L_p[a, b]$. Рассмотрим множество функций $x(t)$, определенных и измеримых на конечном отрезке $[a, b]$ и суммируемых там со степенью p ($p \geq 1$):

$$\int_a^b |x(t)|^p dt < \infty.$$

Множество X является линейным, так как если $x(t) \in X$ и λ — любое число, то $\lambda x(t) \in X$, и если $x(t), y(t) \in X$, сумма их $x(t) + y(t)$ также принадлежит X .*) Норма элемента определяется равенством

$$\|x(t)\|^p = \int_a^b |x(t)|^p dt.$$

Две функции, различающиеся между собой на множестве меры нуль, считаются тождественными.

Выполнение аксиом нормы легко проверяется. То, что $\|x(t)\| \geq 0$ и $\|x(t)\| = 0$ только когда функция $x(t)$ эквивалентна нулю, так же как то, что $\|\lambda x(t)\| = |\lambda| \cdot \|x(t)\|$, является очевидным.

Проверке подлежит только вторая аксиома нормы: $\|x(t) + y(t)\| \leq \|x(t)\| + \|y(t)\|$. Но она является не чем иным, как известным неравенством Минковского

*) Это следует из просто проверяемого неравенства $|\alpha + \beta|^p \leq 2^p (|\alpha|^p + |\beta|^p)$. Для всяких двух чисел α и β верно неравенство $|\alpha + \beta| \leq |\alpha| + |\beta|$. Пусть, для определенности записи, $|\alpha| \geq |\beta|$. Тогда $|\alpha + \beta| \leq 2|\alpha|$, $|\alpha + \beta|^p \leq 2^p |\alpha|^p \leq 2^p (|\alpha|^p + |\beta|^p)$.

$$\left(\int_a^b |x(t) + y(t)|^p dt \right)^{\frac{1}{p}} \leq \left(\int_a^b |x(t)|^p dt \right)^{\frac{1}{p}} + \left(\int_a^b |y(t)|^p dt \right)^{\frac{1}{p}}.$$

Метрика, отвечающая принятой норме, есть

$$\rho(x, y) = \left(\int_a^b |x(t) - y(t)|^p dt \right)^{\frac{1}{p}}.$$

Покажем, что пространство $L_p[a, b]$ является полным. Не ограничивая общности, можно считать отрезок $[a, b]$ совпадающим с $[0, 1]$ и рассматривать пространство $L_p[0, 1]$. Ниже мы будем его обозначать L_p . Возьмем последовательность функций $x_n(t)$ ($n=1, 2, \dots$) из L_p , удовлетворяющую условию Больцано — Коши. Последнее значит, что для всякого $\epsilon_k > 0$ существует такой номер $N(\epsilon_k)$, что

$$\int_0^1 |x_n(t) - x_m(t)|^p dt < \epsilon_k^p$$

для $n, m > N(\epsilon_k)$.

Примем $\epsilon_k^p = \frac{1}{2^{kp+k}}$. При $m, n > N(\epsilon_k)$ неравенство

$$|x_m(t) - x_n(t)| > \frac{1}{2^k}$$

будет выполняться на некотором множестве с мерой, меньшей $\frac{1}{2^k}$.

Если мы возьмем и зафиксируем возрастающую последовательность чисел n_k , причем $n_k > N(\epsilon_k)$, можем сказать, что неравенство

$$|x_{n_{k+1}}(t) - x_{n_k}(t)| > \frac{1}{2^k}$$

выполняется только на некотором множестве A_k с мерой, меньшей $\frac{1}{2^k}$.

Для сокращения обозначим $x_{n_k}(t) = z_k(t)$ и введем множество

$$B_k = [0, 1] - \sum_{i=k}^{\infty} A_i.$$

Очевидно,

$$\text{mes } B_k \geq 1 - \sum_{i=k}^{\infty} \text{mes } A_i > 1 - \left(\frac{1}{2^k} + \frac{1}{2^{k+1}} + \dots \right) = 1 - \frac{1}{2^{k-1}}$$

и, кроме того,

$$|z_{n+1}(t) - z_n(t)| \leq \frac{1}{2^n}$$

для $n > k$ на B_k . Отсюда следует, что $z_n(t)$ на B_k равномерно сходится. Заметив, что $B_k \subset B_{k+1}$, рассмотрим множество

$$B_0 = \lim_{k \rightarrow \infty} B_k = \sum_{k=1}^{\infty} B_k.$$

Из неравенства $\text{mes } B_k > 1 - 2^{-k+1}$ вытекает, что $\text{mes } B_0 = 1$.

Каждая точка B_0 принадлежит некоторому B_k , и последовательность $z_n(t)$ ($n=1, 2, \dots$) поэтому сходится всюду на B_0 . При этом сходимость на каждом B_k является равномерной.

B_0 может отличаться от $[0, 1]$ лишь множеством меры нуль и $z(t)$ определена почти везде на $[0, 1]$.

Последовательность $z_n(t)$ ($n=1, 2, \dots$) есть часть последовательности $x_n(t)$ ($n=1, 2, \dots$), и так как для $x_n(t)$ выполняется признак Больцано — Коши в метрике L_p , то он будет выполняться и для $z_n(t)$:

$$\int_0^1 |z_m(t) - z_n(t)|^p dt < \varepsilon \quad [m, n > n_0(\varepsilon)].$$

Функция, стоящая под знаком интеграла, неотрицательна и верно также неравенство

$$\int_{B_k} |z_m(t) - z_n(t)|^p dt < \varepsilon \quad [m, n > n_0(\varepsilon)].$$

Ввиду равномерной сходимости $z_m(t)$ к $z(t)$ на B_k здесь можно перейти к пределу при $m \rightarrow \infty$ под знаком интеграла

$$\int_{B_k} |z(t) - z_n(t)|^p dt \leq \varepsilon \quad [n > n_0(\varepsilon)].$$

Неравенство выполняется при всяких k . Переходя к пределу при $k \rightarrow \infty$, получим

$$\int_{B_0} |z(t) - z_n(t)|^p dt = \int_0^1 |z(t) - z_n(t)|^p dt \leq \varepsilon \quad [n > n_0(\varepsilon)].$$

Полученное неравенство говорит, в частности, о том, что $z(t) - z_n(t)$ принадлежит L_p , и так как $z_n(t) \in L_p$, то и $z(t) = [z(t) - z_n(t)] + z_n(t)$ также принадлежит L_p .

Кроме того, из неравенства следует, что в метрике L_p $\rho(z, z_n) \leq \frac{1}{n^p}$ $[n > n_0(\varepsilon)]$ и ввиду произвола ε $\rho(z, z_n) \rightarrow 0$ при $n \rightarrow \infty$.

Наконец, $\rho(x_n, z) \leq \rho(x_n, z_k) + \rho(z_k, z) = \rho(x_n, x_{n_k}) + \rho(z_k, z) \rightarrow 0$, когда n и k неограниченно возрастают. Значит, $\rho(x_n, z) \rightarrow 0$ ($n \rightarrow \infty$) и полнота L_p доказана.

Пусть X и Y есть два произвольных множества. Говорят, что на множестве X задан оператор A со значениями из множества Y , если каждому $x \in X$ соответствует некоторый элемент y из Y : $y = A(x)$. Элемент x называют оригиналом и y его изображением. Ниже часто в обозначении оператора $A(x)$ скобки будем опускать и писать Ax .

В частном случае, если Y есть множество чисел и оператор Ax ставит в соответствие каждому элементу $x \in X$ некоторое число, такой оператор A называют функционалом.

Когда множества X и Y являются метрическими пространствами, может быть определено понятие непрерывности оператора. Оператор A называют непрерывным на элементе x , если из сходимости $x_n \rightarrow x$ в пространстве X следует сходимость $Ax_n \rightarrow Ax$ в пространстве Y .

Предположим теперь, что X и Y являются линейными нормированными пространствами. Оператор A называется аддитивным, если для всяких двух элементов x_1 и x_2 из X будет $A(x_1 + x_2) = A(x_1) + A(x_2)$.

Остановим внимание на некоторых свойствах аддитивных операторов.

1. Для аддитивного оператора A верны равенства

$$A(0) = 0, \quad A(-x) = -A(x).$$

В самом деле,

$$A(0) = A(0+0) = A(0) + A(0) \quad \text{и} \quad A(0) = 0.$$

Далее,

$$0 = A(0) = A[x + (-x)] = A(x) + A(-x) \quad \text{и} \quad A(-x) = -A(x).$$

2. Если оператор A аддитивен и непрерывен, то он однороден, т. е. для любого действительного числа λ будет

$$A(\lambda x) = \lambda A(x).$$

Из аддитивности сразу же вытекает, что для всякого целого положительного числа p будет

$$A(px) = pA(x).$$

Если же p есть целое отрицательное число, то

$$A(px) = -pA(-x) = -p[-A(x)] = pA(x).$$

Кроме того, при целом q

$$A(x) = A\left(q \frac{1}{q} x\right) = qA\left(\frac{1}{q} x\right) \quad \text{и} \quad A\left(\frac{1}{q} x\right) = \frac{1}{q} A(x).$$

Поэтому при целых p и q будет

$$A\left(\frac{p}{q} x\right) = \frac{p}{q} A(x).$$

Пусть λ — произвольное действительное число. Возьмем последовательность рациональных чисел $r_n \rightarrow \lambda$. Используя непрерывность A , найдем

$$A(\lambda x) = \lim A(r_n x) = \lim r_n A(x) = \lambda A(x).$$

Аддитивный оператор A , переводящий нормированное пространство X в нормированное пространство Y , называется ограниченным, если существует такое число M , что при всяких $x \in X$ выполняется неравенство

$$\|Ax\| \leq M\|x\|. \quad (1.2)$$

Аддитивный и непрерывный оператор A называют линейным.

Докажем теорему, дающую критерий линейности аддитивного оператора.

Теорема 1. *Для линейности аддитивного оператора A необходимо и достаточно, чтобы он был ограниченным.*

Доказательство. Установим сначала необходимость. Допустим, что оператор A неограничен, и покажем, что это противоречит его непрерывности. Существует последовательность x_n такая, что

$$\|A(x_n)\| \geq n\|x_n\|.$$

Рассмотрим элемент

$$x_n' = \frac{1}{n} \cdot \frac{x_n}{\|x_n\|}.$$

Очевидно, $x_n' \rightarrow 0$ ($n \rightarrow \infty$). Но, с другой стороны,

$$\|Ax_n'\| = \frac{1}{n} \frac{1}{\|x_n\|} \|Ax_n\| \geq 1.$$

Поэтому $A(x_n')$ не стремится к нулевому элементу при $x_n' \rightarrow 0$ и оператор A не является непрерывным на элементе 0.

Достаточность доказывается столь же просто. Если $\|x - x_n\| \rightarrow 0$, то

$$\|A(x) - A(x_n)\| = \|A(x - x_n)\| \leq M\|x - x_n\| \rightarrow 0$$

и, следовательно,

$$A(x_n) \rightarrow A(x).$$

Наименьшее число M , для которого выполняется при любом x неравенство

$$\|A(x)\| \leq M\|x\|,$$

называют нормой оператора A и обозначают $\|A\|$.

Отметим следующее равенство, которое может быть полезным в задаче нахождения нормы:

$$\|A\| = \sup_{\|x\| \leq 1} \|Ax\|. \quad (I.3)$$

Действительно, если $\|x\| \leq 1$, то

$$\|Ax\| \leq \|A\| \cdot \|x\| \leq \|A\|.$$

Поэтому

$$\sup_{\|x\| \leq 1} \|Ax\| \leq \|A\|. \quad (I.3_1)$$

При всяком $\varepsilon > 0$, по определению нормы, существует элемент x' , для которого $\|Ax'\| > (\|A\| - \varepsilon)\|x'\|$. Положим $x = \frac{x'}{\|x'\|}$.

$$\|A(x)\| = \frac{1}{\|x'\|} \|Ax'\| > \frac{1}{\|x'\|} (\|A\| - \varepsilon)\|x'\| = \|A\| - \varepsilon,$$

и так как $\|x\| = 1$, то

$$\sup_{\|x\| \leq 1} \|Ax\| > \|A\| - \varepsilon.$$

Отсюда и из (I.3₁) следует (I.3).

Приведем пример линейного оператора. Рассмотрим преобразование ν -мерного векторного пространства X_ν с элементами $x(x_1, \dots, x_\nu)$ в μ -мерное векторное пространство Y_μ с элементами $y(y_1, \dots, y_\mu)$:

$$y = A(x).$$

При преобразовании любого типа составляющие y_i вектора y будут функциями от x_1, \dots, x_ν :

$$y_i = \varphi_i(x_1, \dots, x_\nu) \quad (i = 1, 2, \dots, \mu).$$

Но если преобразование $y=A(x)$ линейное, φ_i будут аддитивными и непрерывными функциями своих аргументов и, следовательно, будут линейными однородными функциями от x_1, \dots, x_v :

$$y_i = \sum_{j=1}^v a_{ij} x_j \quad (i=1, 2, \dots, \mu). \quad (I.4)$$

В этом случае оператор $A(x)$ есть не что иное, как линейное преобразование с прямоугольной матрицей (a_{ij}) ($i=1, \dots, \mu; j=1, \dots, v$).

Наоборот, преобразование такого вида с любой матрицей (a_{ij}) будет линейным оператором, совершающим преобразование $X_v \rightarrow Y_\mu$.

Рассмотрим вопрос об обращении операторов. Как и выше, будем считать X и Y линейными нормированными пространствами. Пусть A есть аддитивный однородный оператор, отображающий X в Y .

$$A(x) = y. \quad (I.5)$$

Говорят, что оператор A имеет обратный, если существует оператор V , определенный всюду на Y , удовлетворяющий условиям:

1) при всех $x \in X$ верно равенство

$$VA(x) = x;$$

2) при всяких y из Y выполняется равенство

$$AV(y) = y.$$

Оператор V называется двусторонним обратным или просто обратным для A и обозначается $V=A^{-1}$.

Значение каждого из этих условий легко выясняется. Первое из них говорит, что если уравнение (I.5) имеет решение, то оно единственно и имеет представление $x=Vy=A^{-1}y$.

Второе же условие означает, что уравнение (I.5) имеет решение при всяком y и за такое решение может быть принят элемент $x=Vy=A^{-1}y$.

Из определения обратного оператора следует, что $V=A^{-1}$, так же как A , является аддитивным и однородным. В самом деле, если $y_1, y_2 \in Y$, то по второму условию $y_k=A(x_k)$ при $x_k=V(y_k)$ ($k=1, 2$). По первому же условию

$$V(y_1+y_2) = V[A(x_1)+A(x_2)] = VA(x_1+x_2) = x_1+x_2 = V(y_1)+V(y_2).$$

Этим установлена аддитивность V . Подобным образом проверяется однородность V .

Из определения также следует, что оператор A является обратным для A^{-1} , т. е. $(A^{-1})^{-1}=A$.

Покажем теперь, что если оператор A имеет обратный $V=A^{-1}$, то A осуществляет взаимно однозначное отображение X на Y . По второму условию каждый элемент $y \in Y$ есть изображение некоторого элемента $x \in X$. В качестве такого элемента можно взять $x=Vy$. Необходимо лишь проверить, что если $x_1 \neq x_2$, то $A(x_1) \neq A(x_2)$. В самом деле, если бы было $A(x_1)=A(x_2)$, то по первому из условий мы имели бы $x_1=VA(x_1)=VA(x_2)=x_2$.

В некоторых вопросах сходимости вычислительных процессов имеет значение теорема о сходимости последовательности линейных операторов, к доказательству которой мы сейчас перейдем.

Пусть X и Y — два пространства типа В. Рассмотрим последовательность A_n линейных операторов, определенных в X , со значениями из Y .

Последовательность A_n называется сходящейся, если для всякого $x \in X$ в пространстве Y будет сходиться последовательность элементов $y_n=A_nx$. Обозначим $\lim_{n \rightarrow \infty} A_nx = y = Ax$. Оператор A — аддитивный, так

как если в равенстве $A_n(x_1+x_2)=A_nx_1+A_nx_2$ перейти к пределу, то получится $A(x_1+x_2)=Ax_1+Ax_2$. Можно показать, что оператор A будет непрерывным и, стало быть, линейным. Сначала докажем лемму.

Лемма. Если последовательность A_n ($n=1, 2, \dots$) линейных операторов сходится, то нормы их ограничены в совокупности:

$$\|A_n\| \leq M < \infty.$$

Доказательство. Допустим противоположное и будем считать $\sup \|A_n\| = \infty$.

Рассмотрим замкнутый шар радиуса ε с центром в элементе x_0 : $\|x-x_0\| \leq \varepsilon$ и обозначим его $S(x_0, \varepsilon)$. Покажем, что $\|A_nx\|$ не могут быть ограничены в совокупности ни в каком замкнутом шаре.

Пусть, в самом деле, $\|A_nx\| \leq K$ для $x \in S(x_0, \varepsilon)$ ($n=1, 2, \dots$). При всяком $x \in X$ элемент $x' = x_0 + \frac{\varepsilon}{\|x\|}x$ принадлежит $S(x_0, \varepsilon)$, так что

$$\|A_nx'\| = \left\| \frac{\varepsilon}{\|x\|} A_nx + A_nx_0 \right\| \leq K.$$

Отсюда

$$\frac{\varepsilon}{\|x\|} \|A_nx\| - \|A_nx_0\| \leq K$$

и

$$\|A_nx\| \leq \frac{K + \|A_nx_0\|}{\varepsilon} \|x\|.$$

Так как последовательность A_nx_0 сходится и, следовательно, нормы $\|A_nx_0\|$ ограничены в совокупности, то существует такое число K_1 , не за-

висящее от n и x , что $\|A_n x\| \leq K_1 \|x\|$, откуда следует $\|A_n\| \leq K_1$, а это противоречит допущению $\sup \|A_n\| = \infty$.

Возьмем шар $S_0(x_0, \varepsilon_0)$. Ввиду неограниченности в нем $\|A_n x\|$ ($n=1, 2, \dots$) найдется такой оператор A_{n_1} и такой элемент x_1 , что $\|A_{n_1} x_1\| > 1$. Так как оператор A_{n_1} непрерывен, это неравенство будет выполняться в некотором шаре $S_1(x_1, \varepsilon_1)$, содержащемся в $S_0(x_0, \varepsilon_0)$. Для шара S_1 рассуждения могут быть повторены: существует такой оператор A_{n_2} и такой элемент $x_2 \in S_2$, что $\|A_{n_2} x_2\| > 2$ и т. д. Можно считать, что $\varepsilon_n \rightarrow 0$ ($n \rightarrow \infty$). Для построенной последовательности элементов x_1, x_2, \dots будет выполняться признак Больцано — Коши, и она будет сходиться ввиду полноты X к некоторому элементу $x^* \in X$. Элемент x^* будет принадлежать всем шарам S_{n_k} ($k=0, 1, \dots$). При этом $\|A_{n_k} x^*\| > k$. Последнее же противоречит сходимости A_n во всех точках X .

Теперь просто доказывается ограниченность предельного оператора A . Согласно лемме, существует число M такое, что $\|A_n\| \leq M$. Если в неравенстве

$$\|A_n x\| \leq \|A_n\| \|x\| \leq M \|x\|$$

перейти к пределу, получится

$$\|Ax\| \leq M \|x\|$$

и оператор A действительно ограничен. Но так как A аддитивен, он непрерывен и, стало быть, линеен.

Теорема 2 (Банаха—Штейнгауза). *Для сходимости последовательности линейных операторов A_n , отображающих пространство X типа B в пространство Y типа B , необходимо и достаточно выполнение двух условий:*

1) *нормы операторов A_n ограничены в совокупности:*

$$\|A_n\| \leq M < \infty \quad (n=1, 2, \dots);$$

2) *сходимость A_n имеет место на всех элементах множества E , всюду плотного в X .*)*

Доказательство. Необходимость первого условия вытекает из доказанной леммы, необходимость же второго условия очевидна. Остается проверить достаточность условий. Возьмем произвольный элемент $x \in X$ и найдем элемент $\bar{x} \in E$ так, чтобы было $\|x - \bar{x}\| \leq \frac{\varepsilon}{3M}$. Последовательность $A_n \bar{x}$ сходится, и для больших n будет верно неравенство

$$\|A_{n+m} \bar{x} - A_n \bar{x}\| \leq \frac{\varepsilon}{3}.$$

*) Множество E называется всюду плотным в X , если каждый элемент $x \in X$ может быть приближен по норме сколь угодно точно элементом из E .

Для таких n будет

$$\begin{aligned} \|A_{n+m}x - A_nx\| &\leq \|A_{n+m}x - A_{n+m}\bar{x}\| + \|A_{n+m}\bar{x} - A_n\bar{x}\| + \\ &+ \|A_n\bar{x} - A_nx\| \leq 2M\|x - \bar{x}\| + \frac{\varepsilon}{3} < \frac{2\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon. \end{aligned}$$

Для последовательности A_nx признак Больцано — Коши выполняется, ввиду же полноты пространства Y , последовательность будет сходящейся:

$$\lim_{n \rightarrow \infty} A_nx = y = Ax$$

при всяком $x \in X$.

Для некоторых приложений полезна видоизмененная теорема Банаха — Штейнгауза, в которой устанавливаются условия сходимости последовательности операторов к заданному оператору.*)

Как и выше, предположим, что X и Y есть пространства типа В и допустим, что в X определены линейный оператор A^* и последовательность линейных операторов A_n ($n=1, 2, \dots$) со значениями в Y .

Теорема 2' (Банаха—Штейнгауза). *Для сходимости последовательности операторов A_n к A^* необходимо и достаточно выполнение условий:*

1') *нормы операторов A_n ограничены в совокупности*

$$\|A_n\| \leq M < \infty \quad (n=1, 2, \dots);$$

2') *на всех элементах x множества E , всюду плотного в X , имеет место сходимость*

$$A_nx \rightarrow A^*x.$$

Эта теорема является очевидным следствием теоремы 2. Проверке подлежит только достаточность условий 1' и 2'. Если они выполняются, то, по теореме 2, существует линейный оператор A такой, что для всякого x из X , будет $A_nx \rightarrow Ax$. По второму же условию на всех элементах множества E должно быть $A=A^*$. Но если два линейных оператора совпадают на всюду плотном в X множестве, то они равны.

*) Например, когда рассматривается процесс вычисления определенного интеграла при помощи приближенных квадратурных правил

$$I(f) = \int_a^b p(x)f(x)dx \approx \sum_{k=1}^n A_k^{(n)} f[x_k^{(n)}] = Q_n(f),$$

то вопрос о сходимости этого процесса с точки зрения функционального анализа есть проблема сходимости последовательности операторов Q_n к интегральному оператору I . Аналогичное можно сказать о процессе интерполирования функции, о процессе разложения функции в ряд Фурье и др.

В самом деле, пусть x и x' есть два произвольных элемента соответственно из X и E . Очевидно,

$$Ax - A^*x = [A(x - x') + Ax'] - [A^*(x - x') + A^*x'] = A(x - x') - A^*(x - x').$$

Поэтому

$$\|Ax - A^*x\| \leq (\|A\| + \|A^*\|)\|x - x'\|.$$

Левая часть неравенства не зависит от x' . Правая же часть может быть сделана меньше любого числа, так как путем выбора x' множитель $\|x - x'\|$ может быть сделан сколь угодно малым. Отсюда следует, что

$$\|Ax - A^*x\| = 0, \quad Ax = A^*x$$

и оператор A совпадает с A^* всюду в X .

§ 3. ДИФФЕРЕНЦИРОВАНИЕ НЕЛИНЕЙНЫХ ОПЕРАТОРОВ И НЕКОТОРЫЕ ТЕОРЕМЫ, С ЭТИМ СВЯЗАННЫЕ

Предварительно ознакомимся с некоторыми фактами теории линейных и билинейных операторов, необходимыми для изложения теории дифференцирования.

Пусть X и Y два линейных нормированных пространства. Рассмотрим множество линейных операторов, переводящих X в Y . Такое множество обозначим символом $[X \rightarrow Y]$. Убедимся в том, что при надлежащем определении действий сложения, умножения на число и выбора нормы его можно сделать линейным нормированным пространством.

Возьмем два произвольных элемента A_1 и A_2 из $[X \rightarrow Y]$. Под суммой их $A = A_1 + A_2$ условимся понимать оператор из X в Y , определенный равенством

$$A(x) = A_1(x) + A_2(x).$$

A есть, очевидно, аддитивный оператор. Кроме того, ввиду

$$\|A(x)\| \leq \|A_1(x) + A_2(x)\| \leq (\|A_1\| + \|A_2\|)\|x\|,$$

A есть ограниченный, следовательно, линейный оператор и $A \in [X \rightarrow Y]$. Из указанного неравенства получается оценка его нормы:

$$\|A\| \leq \|A_1\| + \|A_2\|. \quad (I.6)$$

Далее, пусть λ есть произвольное число и A — линейный оператор из $[X \rightarrow Y]$. Под произведением $\lambda A = \tilde{A}$ понимается оператор, определенный правилом

$$\tilde{A}(x) = \lambda A(x).$$

Как и выше, легко убедимся в том, что $\tilde{A} \in [X \rightarrow Y]$ и

$$\|\tilde{A}\| = |\lambda| \cdot \|A\|. \quad (I.7)$$

В качестве нулевого элемента в множестве $[X \rightarrow Y]$ выберем оператор A_0 , значения которого тождественно равны нулю:

$$A_0(x) = 0 \text{ при } x \in X. \quad (I.8)$$

Можно проверить, что все аксиомы линейного множества, указанные в начале § 2, будут выполняться.

Нам осталось еще определить норму элементов в $[X \rightarrow Y]$.

За норму A как элемента $[X \rightarrow Y]$ примем норму оператора $A(x)$:

$$\|A\| = \sup_{\|x\| \leq 1} \|A(x)\|.$$

Выполнение аксиом нормы, приведенных в начале § 2, легко проверить, опираясь на (I.6) — (I.8).

Введем понятие о билинейном операторе.

Пусть $B(x', x)$ каждой паре элементов $x', x \in X$ ставит в соответствие элемент $y \in Y$:

$$y = B(x', x).$$

$B(x', x)$ называется билинейным оператором, если выполняются условия:

1) для $B(x', x)$ верны равенства

$$\left. \begin{aligned} B(\alpha x' + \tilde{\alpha} \tilde{x}', x) &= \alpha B(x', x) + \tilde{\alpha} B(\tilde{x}', x), \\ B(x', \alpha x + \tilde{\alpha} \tilde{x}) &= \alpha B(x', x) + \tilde{\alpha} B(x', \tilde{x}) \end{aligned} \right\} \quad (I.9)$$

при любых числах $\alpha, \tilde{\alpha}$ и элементах $x, x', \tilde{x}, \tilde{x}'$;

2) существует число M такое, что при всяких $x', x \in X$ выполняется неравенство

$$\|B(x', x)\| \leq M \|x'\| \|x\|. \quad (I.10)$$

Наименьшее возможное значение M называется нормой оператора $B(x', x)$ и обозначается $\|B\|$.

Приведем простой пример билинейного оператора. Допустим, что X_ν есть ν -мерное векторное пространство с элементами $x(x_1, x_2, \dots, x_\nu)$ и Y_μ есть μ -мерное векторное пространство с элементами $y(y_1, y_2, \dots, y_\mu)$.

Рассмотрим билинейный оператор $y = B(x', x)$ из X_ν в Y_μ . Каждая составляющая y_i элемента y будет, очевидно, билинейной формой составляющих $x'(x'_1, \dots, x'_\nu)$ и $x(x_1, \dots, x_\nu)$:

$$y_i = \sum_{j, k=1}^{\nu} a_{jk}^{(i)} x_j' x_k \quad (i=1, 2, \dots, \mu). \quad (\text{I.11})$$

Поэтому $y=B(x', x)$ будет μ -мерный вектор, составляющие которого есть билинейные формы, определенные последним равенством.

Что касается нормы оператора B , то ее численное значение зависит от того, как определена норма в X_{ν} и Y_{μ} . Например, будем считать, что в X_{ν} и Y_{μ} взяты кубические нормы вектора:

$$\|x\|_m = \max_i |x_i| \quad \text{и} \quad \|y\|_m = \max_i |y_i|.$$

$$\begin{aligned} \|y\|_m &= \|B(x', x)\|_m = \max_i \left| \sum_{j, k=1}^{\nu} a_{jk}^{(i)} x_j' x_k \right| \leq \\ &\leq \max_i \sum_{j, k=1}^{\nu} |a_{jk}^{(i)}| \cdot \|x'\|_m \cdot \|x\|_m \end{aligned}$$

и

$$\|B\|_m \leq \max_i \sum_{j, k=1}^{\nu} |a_{jk}^{(i)}|.$$

Как мы увидим несколькими страницами далее, при определении второй производной от оператора нам придется иметь дело с линейными отображениями пространства X на пространство линейных операций $[X \rightarrow Y]$. Множество таких операторов кратко будем обозначать $[X \rightarrow [X \rightarrow Y]]$. Покажем сейчас, что между этими отображениями и билинейными операторами $B(x', x)$ существует простая связь, позволяющая сказать, что, по существу дела, равносильно рассматривать $B(x', x)$ как билинейный оператор, или как оператор отображения $X \rightarrow [X \rightarrow Y]$. Это обстоятельство позволит упростить изложение, так как билинейные операторы B являются более наглядным аппаратом, чем операторы из $[X \rightarrow [X \rightarrow Y]]$, и часто более удобны для вычислений.

Пусть U есть один из операторов, отображающих X в $[X \rightarrow Y]$. Возьмем произвольный элемент $x' \in X$ и положим $A_{x'} = U(x')$.

Оператор $A_{x'}$ принадлежит $[X \rightarrow Y]$ и $y = A_{x'}(x) = U(x')x = B(x', x)$ является элементом Y . Построенный так оператор $B(x', x)$ будет удовлетворять второму условию (I.9) по аргументу x , так как $A_{x'}(x)$ есть линейный оператор относительно x , и будет удовлетворять первому из условий (I.9) по x' , так как $U(x')$ линеен относительно x' .

Условие (I.10) также выполняется, так как

$$\|B(x', x)\| \leq \|A_{x'}\| \cdot \|x\| = \|U(x')\| \cdot \|x\| \leq \|U\| \cdot \|x'\| \cdot \|x\|.$$

Поэтому $B(x', x)$ есть билинейный оператор, при этом

$$\|B\| \leq \|U\|. \quad (\text{I.12})$$

Таким образом, каждому оператору U из $[X \rightarrow [X \rightarrow Y]]$ указанным способом ставится в соответствие билинейный оператор $B(x', x)$. Проверим теперь, что каждый билинейный оператор $B(x', x)$ при указанном законе соответствия является образом некоторого оператора U из $[X \rightarrow [X \rightarrow Y]]$. В самом деле, если мы фиксируем x' , то оператор $U(x') = B(x', \cdot)$ будет принадлежать множеству $[X \rightarrow Y]$, и так как

$$\|U(x')\| = \sup_{\|x\| \leq 1} \|U(x')(x)\| = \sup_{\|x\| \leq 1} \|B(x', x)\| \leq \|B\| \cdot \|x'\|,$$

то оператор U линеен и принадлежит множеству $[X \rightarrow [X \rightarrow Y]]$. Кроме того, $\|U\| \leq \|B\|$. Сравнение последнего неравенства с (I.12) приводит к заключению:

$$\|U\| = \|B\|. \quad (\text{I.13})$$

Дадим определение производной первого порядка от оператора. Будем, как выше, считать X и Y линейными нормированными пространствами. Рассмотрим оператор

$$y = f(x),$$

переводящий X в Y . Говорят, что оператор f дифференцируем (по Фреше) на элементе x , если существует линейный оператор $H \in [X \rightarrow Y]$ такой, что

$$\|f(x + \Delta x) - f(x) - H(\Delta x)\| \leq \|\Delta x\| \cdot \varepsilon(\|\Delta x\|), \quad (\text{I.14})$$

где $\varepsilon(\delta) \rightarrow 0$ при $\delta \rightarrow 0$. Оператор H называют производной оператора f на элементе x :

$$H = f'(x).$$

Линейное преобразование $H(\Delta x)$ имеет смысл дифференциала оператора.

Рассмотрим как пример оператор, преобразующий ν -мерное векторное пространство X_ν в μ -мерное векторное пространство Y_μ . Такое преобразование $y = f(x)$ определяется совокупностью μ функций

$$y_i = f_i(x_1, x_2, \dots, x_\nu) \quad (i = 1, 2, \dots, \mu). \quad (\text{I.15})$$

Функции f_i будем предполагать дифференцируемыми.

По определению производной (I.14), мы должны из изменения оператора $f(x + \Delta x) - f(x)$ выделить такой линейный оператор $H(\Delta x)$, аргументом которого является изменение Δx вектора x [$\Delta x = (\Delta x_1, \Delta x_2, \dots, \Delta x_\nu)$], что

$$f(x + \Delta x) - f(x) - H(\Delta x) = \varepsilon(x, \Delta x),$$

где

$$\frac{\|\varepsilon(x, \Delta x)\|}{\|\Delta x\|} \rightarrow 0 \quad (\|\Delta x\| \rightarrow 0).$$

^{*} Под $B(x', \cdot)$ понимается такой оператор V , что $V(x) = B(x', x)$.

В предыдущем параграфе мы выяснили, какую форму имеет линейный оператор из X_v в Y_μ , и можно сказать, что преобразование $H(\Delta x)$ будет иметь вид

$$\{H(\Delta x)\}_i = \sum_{j=1}^v a_{ij} \Delta x_j \quad (i=1, 2, \dots, \mu). \quad (I.16)$$

Мы должны определить для него коэффициенты a_{ij} .

Изменение оператора $f(x+\Delta x) - f(x)$ будет вектором с μ составляющими

$$f_i(x_1 + \Delta x_1, \dots, x_v + \Delta x_v) - f_i(x_1, \dots, x_v) = f_i(x + \Delta x) - f_i(x) = \Delta f_i(x). \quad (I.17)$$

Из приращений функций $\Delta f_i(x)$ нам следует выделить, согласно с (I.16), главную часть, линейно зависящую от $\Delta x_1, \dots, \Delta x_v$. Но такая главная часть есть дифференциал функции f_i в точке (x_1, \dots, x_v)

$$\begin{aligned} \{H(\Delta x)\}_i &= \sum_{j=1}^v \frac{\partial f_i(x)}{\partial x_j} \Delta x_j \quad (i=1, 2, \dots, \mu), \\ a_{ij} &= \frac{\partial f_i}{\partial x_j}. \end{aligned}$$

Таким образом, производная $f'(x)$ есть линейный оператор, осуществляющий линейное преобразование $X_v \rightarrow Y_\mu$, определяемое матрицей Якоби

$$f'(x) = \left[\frac{\partial f_i}{\partial x_j} \right] \quad \left(\begin{array}{l} i=1, 2, \dots, \mu \\ j=1, 2, \dots, v \end{array} \right).$$

Перейдем к определению второй производной от f . Первая производная $f'(x)$ есть линейный оператор, переводящий X в Y , и, стало быть, он является элементом пространства линейных операторов $(X \rightarrow Y)$. Этот оператор зависит от элемента x , как от параметра, и для каждого x будет своим, подобно тому как производная $F'(t)$ от функции $F(t)$ зависит от положения точки t , в которой она вычисляется. $H = f'(x)$ можно рассматривать как оператор, преобразующий пространство X в пространство $[X \rightarrow Y]$. Может оказаться, что этот последний оператор будет дифференцируемым. Производная от $f'(x)$ называется второй производной от f и обозначается $f''(x)$:

$$V = [f'(x)]' = f''(x).$$

$f''(x)$ есть элемент пространства $[X \rightarrow [X \rightarrow Y]]$, т. е. это линейный оператор, преобразующий X в $[X \rightarrow Y]$. В начале настоящего параграфа мы

выяснили, что рассмотрение такого оператора равносильно рассмотрению билинейного оператора $y = B(x', x)$ из X в Y . В соответствии с этим под $\|f''(x)\|$ мы будем понимать норму билинейного оператора.

Выполним вычисление второй производной $f''(x)$ для частного случая, когда $f(x)$ есть рассмотренный несколькими строками выше оператор, преобразующий ν -мерное векторное пространство X_ν в μ -мерное векторное пространство Y_μ . Функции $f_i(x_1, \dots, x_\nu)$ в (I.15) будем считать дважды дифференцируемыми.

Предположим, что вторая производная $f''(x)$ существует и найдем лишь ее значение. Возьмем произвольный элемент $x_0 \in X_\nu$ и вычислим вторую производную $f''(x_0)$ на x_0 .

По определению (I.14), для вычисления $f''(x_0)$ нужно найти оператор $U \in [X \rightarrow [X \rightarrow Y]]$ такой, для которого выполняется условие

$$\|f'(x_0 + \Delta x) - f'(x_0) - U(\Delta x)\| \leq \|\Delta x\| \varepsilon(\|\Delta x\|). \quad (\text{I.18})$$

Мы преобразуем это условие к другой форме, более удобной для вычислений. Полагая $\Delta x = tx'$ и разделив обе части (I.18) на t , мы, как следствие, получим соотношение

$$\lim_{t \rightarrow 0} \frac{f'(x_0 + tx') - f'(x_0)}{t} = U(x') = A_{x'}. \quad (\text{I.19})$$

В левой и правой частях здесь стоят линейные операторы, принадлежащие $[X \rightarrow [X \rightarrow Y]]$. Нам удобнее найти вторую производную не в форме оператора U , а при помощи соответствующей ему билинейной формы $B(x', x)$. Для этого возьмем произвольный элемент $x \in X$ и, заметив, что $A_{x'}(x) = U(x')(x) = B(x', x)$, от (I.19) перейдем к равенству

$$\lim_{t \rightarrow 0} \frac{1}{t} [f'(x_0 + tx')x - f'(x_0)x] = B(x', x). \quad (\text{I.20})$$

Значение билинейного оператора $B(x', x)$ принадлежит пространству Y_μ и является вектором $y(y_1, \dots, y_\mu)$, составляющие y_i которого являются билинейными формами от (x_1, \dots, x_ν) , (x'_1, \dots, x'_ν) и даны равенством (I.11)

$$\{B(x', x)\}_i = \sum_{j, k=1}^{\nu} a_{jkh}^{(i)} x'_j x_k. \quad (\text{I.21})$$

Нашей задачей является вычисление $a_{jkh}^{(i)}$. Значение $f'(x_0)x$ также является вектором из Y , и составляющие его были вычислены раньше:

$$\{f'(x_0)x\}_i = \sum_{k=1}^{\nu} \frac{\partial f_i(x_0)}{\partial x_k} x_k.$$

Аналогичные выражения с заменой x_0 на $x_0 + tx'$ верны для составляющих $f'(x_0 + tx')x$.

Векторное равенство (I.20) равносильно следующим μ численным равенствам:

$$\lim_{t \rightarrow 0} \frac{1}{t} \left[\sum_{k=1}^{\nu} \left(\frac{\partial f_i(x_0 + tx')}{\partial x_k} - \frac{\partial f_i(x_0)}{\partial x_k} \right) x_k \right] = \{B(x', x)\}_i.$$

Отсюда сразу же получается

$$\sum_{j, k=1}^{\nu} \frac{\partial^2 f_i(x_0)}{\partial x_j \partial x_k} x_j' x_k = \{B(x', x)\}_i \quad (i=1, 2, \dots, \mu)$$

и

$$a_{jk}^{(i)} = \frac{\partial^2 f_i(x_0)}{\partial x_j \partial x_k}. \quad (I.22)$$

Докажем теперь две теоремы, которые используются при изложении метода Ньютона для решения операторных уравнений.

Сначала приведем два простых утверждения, справедливость которых просто проверяется на основании определения (I.14) первой производной.

Лемма 1. Пусть f есть линейный оператор из X в Y и x есть произвольный элемент из X .

Тогда

$$f'(x) = f.$$

Лемма 2. Пусть X, Y, Z — линейные нормированные пространства. Если H есть линейный оператор из Y в Z и $y = f(x)$ дифференцируемый оператор из X в Y , то

$$[Hf(x)]' = Hf'(x).$$

Ниже сформулирована и доказывается теорема об изменении оператора.

Теорема 1. Если $f(x)$ есть дифференцируемый оператор, то верно неравенство

$$\|f(x + \Delta x) - f(x)\| \leq \sup_{0 \leq \theta \leq 1} \|f'(x + \theta \Delta x)\| \cdot \|\Delta x\|. \quad (I.23)$$

Доказательство. Положим $f(x + \Delta x) - f(x) = y$ и рассмотрим в пространстве Y линейный функционал T , обладающий свойствами

$$\|T\| = 1, \quad T(y) = \|y\|.$$

Существование такого функционала может быть доказано, и на этом мы останавливаться не будем.*) Образует функцию действительного аргумента t : $F(t) = T[f(x+t\Delta x)]$. Ее производная вычисляется при помощи леммы 2:

$$F'(t) = T \cdot f'(x+t\Delta x) \Delta x.$$

На основании известной в анализе теоремы о конечном приращении

$$T(y) = T[f(x+\Delta x) - f(x)] = F(1) - F(0) = F'(\theta) = (Tf'(x+\theta\Delta x)) \Delta x.$$

Отсюда получается сразу же доказательство теоремы:

$$\begin{aligned} \|f(x+\Delta x) - f(x)\| &= T(y) \leq \|T\| \cdot \|f'(x+\theta\Delta x)\| \cdot \|\Delta x\| \leq \\ &\leq \max_{0 \leq \theta \leq 1} \|f'(x+\theta\Delta x)\| \cdot \|\Delta x\|. \end{aligned}$$

Теорема 2. Если $f(x)$ есть дважды дифференцируемый оператор, то справедлива оценка

$$\|f(x+\Delta x) - f(x) - f'(x)\Delta x\| \leq \frac{1}{2} \max_{0 \leq \theta \leq 1} \|f''(x+\theta\Delta x)\| \cdot \|\Delta x\|^2.$$

Доказательство. В доказательстве будут в значительной мере повторены проделанные в предыдущем случае рассуждения. Обозначим

$$y = f(x+\Delta x) - f(x) - f'(x)\Delta x$$

и построим вспомогательную функцию численного аргумента t :

$$F(t) = T(f(x+t\Delta x)),$$

где T — тот же функционал, что и в доказательстве предыдущей теоремы. Для $F(t)$ находим следующие значения производных:

$$F'(t) = T[f'(x+t\Delta x)\Delta x], \quad F''(t) = T[f''(x+t\Delta x)\Delta x\Delta x],$$

при этом последняя запись означает, что билинейный оператор $f''(x+t\Delta x)$ должен быть вычислен для одинаковых значений аргументов, равных Δx . После этого остается лишь воспользоваться формулой Тейлора с остатком:

$$\begin{aligned} \|y\| &= T(y) = F(1) - F(0) - F'(0) = \\ &= \frac{1}{2} F''(\theta) \leq \frac{1}{2} \sup_{0 \leq \theta \leq 1} \|f''(x+\theta\Delta x)\| \cdot \|\Delta x\|^2. \end{aligned}$$

*) Это есть простое следствие теоремы Банаха — Хана о продолжении линейного функционала. См., например, Л. А. Люстерник и В. И. Соболев. Элементы функционального анализа, гл. III, § 21. М., 1951.

Добавление II

ЧИСЛА И МНОГОЧЛЕНЫ БЕРНУЛЛИ

§ 1. ЧИСЛА БЕРНУЛЛИ

Определим их при помощи производящей функции. Пусть t есть комплексная переменная. Рассмотрим функцию

$$g(t) = \frac{t}{e^t - 1}. \quad (\text{II.1})$$

Она регулярна в круге $|t| < 2\pi$ и может быть там разложена в степенной ряд по t . Запишем разложение в форме

$$\frac{t}{e^t - 1} = \sum_{n=0}^{\infty} \frac{B_n}{n!} t^n. \quad (\text{II.2})$$

Определенные этим равенством величины B_n и называются числами Бернулли.

Можно легко построить рекурсионное соотношение, позволяющее последовательно вычислять B_n . Умножим обе части (II.2) на $e^t - 1 =$

$$= \sum_{v=1}^{\infty} \frac{t^v}{v!}.$$

$$\left(\frac{t}{1!} + \frac{t^2}{2!} + \frac{t^3}{3!} + \dots \right) \sum_{n=0}^{\infty} \frac{B_n}{n!} t^n = t.$$

Сравнение коэффициентов при t, t^2, t^3, \dots дает нужные нам соотношения:

$$B_0 = 1, \quad \frac{B_0}{n!} + \frac{B_1}{(n-1)!1!} + \frac{B_2}{(n-2)!2!} + \dots + \frac{B_{n-1}}{1!(n-1)!} = 0 \quad (\text{II.3})$$
$$(n=2, 3, \dots).$$

Последнему равенству можно придать удобную для запоминания форму.

Умножив обе части равенства на $n!$ и прибавив к ним B_n , получим

$$\sum_{k=0}^n \frac{n!}{k!(n-k)!} B_k = B_n.$$

Левая часть равенства аналогична степени двучлена, и само равенство можно записать в условном виде

$$(1+B)^n = B_n, \quad (\text{II.4})$$

где после возведения двучлена в степень нужно показатели степеней B толковать как индексы чисел Бернулли. Покажем, что все числа Бернулли с нечетными индексами, большими единицы, равны нулю:

$$B_{2k+1} = 0 \quad (k=1, 2, \dots). \quad (\text{II.5})$$

Заменим в (II.2) t на $-t$:

$$\frac{-t}{e^{-t}-1} = \sum_{n=0}^{\infty} (-1)^n \frac{B_n}{n!} t^n.$$

Но

$$\frac{-t}{e^{-t}-1} = \frac{et \cdot t}{e^t-1} = t + \frac{t}{e^t-1} = t + \sum_{n=0}^{\infty} \frac{B_n}{n!} t^n$$

и, следовательно, должно быть

$$t + \sum_{n=0}^{\infty} \frac{B_n}{n!} t^n = \sum_{n=0}^{\infty} (-1)^n \frac{B_n}{n!} t^n.$$

Сравнивая здесь коэффициенты при t^n ($n > 1$), получим $B_n = (-1)^n B_n$. Для $n = 2k+1$ это дает $B_{2k+1} = -B_{2k+1}$ и $B_{2k+1} = 0$ ($k=1, 2, \dots$).

Приведем значения нескольких первых чисел Бернулли:

$$B_0 = 1, B_1 = -\frac{1}{2}, B_2 = \frac{1}{6}, B_4 = -\frac{1}{30}, B_6 = \frac{1}{42}, B_8 = -\frac{1}{30},$$

$$B_{10} = \frac{5}{66}, B_{12} = -\frac{691}{2730}, B_{14} = \frac{7}{6}, B_{16} = -\frac{3617}{510}, \dots$$

Чтобы закончить перечисление нужных свойств чисел Бернулли, укажем еще на связь их с обратными степенями целых чисел:

$$B_{2k} = \frac{(-1)^{k-1}(2k)!}{2^{2k-1}\pi^{2k}} \left(1 + \frac{1}{2^{2k}} + \frac{1}{3^{2k}} + \frac{1}{4^{2k}} + \dots \right). \quad (\text{II.6})$$

Это равенство следует сразу же из доказанного ниже разложения (II.18) многочлена Бернулли в тригонометрический ряд на $[0, 1]$.

Полезно отметить два следствия, вытекающих из (II.6): знак B_{2k} совпадает с $(-1)^{k-1}$ и числа Бернулли смежных четных номеров всегда противоположны по знаку.

При больших k погрешность приводимого ниже приближенного равенства будет малой величиной:

$$B_{2k} \approx 2(-1)^{k-1}(2k)!(2\pi)^{-2k}.$$

Из него следует, что с увеличением k число B_{2k} будет быстро возрастать.

§ 2. МНОГОЧЛЕНЫ БЕРНУЛЛИ И ИХ СВОЙСТВА

Возьмем функцию

$$g(x, t) = e^{xt} \frac{t}{e^t - 1}, \quad (\text{II.7})$$

отличающуюся от (II.1) множителем e^{xt} . Она регулярна в круге $|t| < 2\pi$ и может быть разложена там в степенной ряд

$$g(x, t) = e^{xt} \frac{t}{e^t - 1} = \sum_{n=0}^{\infty} \frac{B_n(x)}{n!} t^n. \quad (\text{II.8})$$

Ниже мы увидим, что коэффициент $B_n(x)$ является многочленом степени n . Он называется многочленом Бернулли. Найдем явное выражение его через числа Бернулли и степени x .

Заменим в (II.8) e^{xt} рядом $\sum_{v=0}^{\infty} \frac{x^v t^v}{v!}$ и $\frac{t}{e^t - 1}$ — разложением (II.2):

$$\sum_{v=0}^{\infty} \frac{x^v t^v}{v!} \sum_{n=0}^{\infty} \frac{B_n}{n!} t^n = \sum_{n=0}^{\infty} \frac{B_n(x)}{n!} t^n.$$

Сравнив коэффициенты при t^n , получим равенство

$$\frac{B_n(x)}{n!} = \frac{x^n B_0}{n!} + \frac{x^{n-1} B_1}{(n-1)!1!} + \dots + \frac{B_n}{n!}$$

или после умножения на $n!$

$$B_n(x) = \sum_{k=0}^n \frac{n!}{k!(n-k)!} B_{n-k} x^k, \quad (\text{II.9})$$

что можно записать в простой условной форме:

$$B_n(x) = (x+B)^n.$$

Ознакомимся с некоторыми свойствами многочленов Бернулли.

1. Начальное значение многочлена Бернулли при $x=0$ равно числу Бернулли того же номера:

$$B_n(0) = B_n. \quad (\text{II.10})$$

2. Дифференцирование и интегрирование $B_n(x)$. Вычислив производную по x от обеих частей равенства (II.8), получим

$$te^{xt} \frac{t}{e^t - 1} = \sum_{n=0}^{\infty} \frac{B_n'(x)}{n!} t^n.$$

С другой стороны,

$$te^{xt} \frac{t}{e^t - 1} = \sum_{n=0}^{\infty} \frac{B_n(x)}{n!} t^{n+1}.$$

Левые части обоих последних равенств одинаковы, сравнение же коэффициентов разложений, стоящих в правых частях, дает

$$\frac{B_n'(x)}{n!} = \frac{B_{n-1}(x)}{(n-1)!}$$

или

$$B_n'(x) = nB_{n-1}(x). \quad (\text{II.11})$$

Отсюда, если воспользоваться (II.10), получим правило интегрирования

$$B_n(x) = B_n + n \int_0^x B_{n-1}(t) dt. \quad (\text{II.12})$$

3. Симметрия распределения значений $B_n(x)$. Точки x и $1-x$ расположены симметрично относительно точки $x = \frac{1}{2}$, являю-

шейся серединой отрезка $[0, 1]$. Мы покажем, что при всяком n и любых x выполняется равенство $B_n(1-x) = (-1)^n B_n(x)$. В частности, отсюда следует, что $B_{2k}(1-x) = B_{2k}(x)$ и график многочлена Бернулли четного индекса $2k$ симметричен относительно прямой $x = \frac{1}{2}$. При $n = 2k+1$ будет $B_{2k+1}(1-x) = -B_{2k+1}(x)$ и график $B_{2k+1}(x)$ антисимметричен относительно прямой $x = \frac{1}{2}$.

Для доказательства заменим в (II.8) x на $1-x$:

$$\begin{aligned} \sum_{n=0}^{\infty} \frac{B_n(1-x)}{n!} t^n &= e^{(1-x)t} \frac{t}{e^t - 1} = e^{-xt} \frac{e^t t}{e^t - 1} = \\ &= e^{-xt} \frac{-t}{e^{-t} - 1} = \sum_{n=0}^{\infty} \frac{B_n(x)}{n!} (-t)^n. \end{aligned}$$

Сравнение коэффициентов при t^n в первом и последнем разложениях приводит к нужному равенству

$$B_n(1-x) = (-1)^n B_n(x). \quad (\text{II.13})$$

4. Изменение $B_n(x)$ на отрезке $[0, 1]$. Будем рассматривать многочлены $y_n(x) = B_n(x) - B_n$, несущественно отличающиеся от $B_n(x)$. Пусть $n > 1$. Покажем, что точки $x=0$ и $x=1$ являются нулями $y_n(x)$. Действительно, на основании (II.5), (II.10) и (II.13)

$$y_n(0) = B_n(0) - B_n = B_n - B_n = 0$$

и

$$y_n(1) = B_n(1) - B_n = (-1)^n B_n(0) - B_n = -B_n[1 - (-1)^n] = 0,$$

ввиду того что при четном n равна нулю величина, стоящая в прямых скобках, а при нечетном n , большем единицы, равно нулю число Бернулли B_n .

Рассмотрим теперь многочлен нечетного номера $y_{2k+1}(x) = B_{2k+1}(x)$ ($k > 0$). Значение $x = \frac{1}{2}$ есть нуль $y_{2k+1}(x)$, так как из (II.13) следует

$$B_{2k+1}\left(\frac{1}{2}\right) = -B_{2k+1}\left(\frac{1}{2}\right)$$

и, значит,

$$B_{2k+1}\left(\frac{1}{2}\right) = y_{2k+1}\left(\frac{1}{2}\right) = 0.$$

Убедимся теперь, что внутри $[0, 1]$ $y_{2k+1}(x)$ не имеет корней, отличных от $x = \frac{1}{2}$. Для этого достаточно показать, что $y_{2k+1}(x)$ не может иметь внутри $[0, 1]$ двух разных нулей. Допустим противоположное: пусть α и β ($0 < \alpha < \beta < 1$) являются нулями $y_{2k+1}(x)$. Ввиду того что $x=0$ и $x=1$ также являются нулями, внутри каждого из отрезков $[0, \alpha]$, $[\alpha, \beta]$ и $[\beta, 1]$ многочлен

$$y'_{2k+1}(x) = B'_{2k+1}(x) = (2k+1)B_{2k}(x)$$

должен иметь по меньшей мере один нуль и, следовательно, многочлен

$$y''_{2k+1}(x) = (2k+1)B'_{2k}(x) = (2k+1)2ky_{2k-1}(x),$$

а стало быть и $y_{2k-1}(x)$, должен иметь внутри $[0, 1]$ по меньшей мере два разных нуля. Продолжая рассуждения, мы убедимся в том, что многочлен $y_3(x)$ имеет внутри $[0, 1]$ не меньше двух разных корней. Если к ним присоединить два корня $x=0$ и $x=1$, мы придем к невозможному заключению, что многочлен третьей степени, отличный от тождественного нуля, $y_3(x)$ имеет не меньше четырех разных нулей. Поэтому наше допущение о том, что $y_{2k+1}(x)$ имеет два разных нуля между точками $x=0$ и $x=1$ является неверным.

Обратимся теперь к многочленам четного индекса $y_{2k}(x)$. Из только что доказанного вытекает, что внутри $[0, 1]$ $y_{2k}(x)$ не может иметь корней и, следовательно, сохраняет знак на $[0, 1]$. В самом деле, если бы $y_{2k}(x)$ обращался в нуль внутри $[0, 1]$, производная

$$y'_{2k}(x) = B'_{2k}(x) = 2kB_{2k-1}(x) = 2ky_{2k-1}(x),$$

а следовательно, и многочлен $y_{2k-1}(x)$ должны были бы иметь внутри $[0, 1]$ не меньше двух разных корней, что невозможно.

Знак $y_{2k}(x)$ на $[0, 1]$ можно определить, если вычислить его значение в одной внутренней точке $[0, 1]$, например в точке $x = \frac{1}{2}$. Положив в (II.8) $x = \frac{1}{2}$ и выполнив несложные преобразования, получим цепочку равенств

$$\begin{aligned} \sum_{n=0}^{\infty} \frac{B_n \left(\frac{1}{2} \right)}{n!} t^n &= e^{\frac{1}{2}t} \frac{t}{e^t - 1} = \frac{(e^{\frac{1}{2}t} + 1)t}{e^t - 1} - \frac{t}{e^t - 1} = \\ &= 2 \frac{\frac{1}{2}t}{e^{\frac{1}{2}t} - 1} - \frac{t}{e^t - 1} = \end{aligned}$$

$$= 2 \sum_{n=0}^{\infty} \frac{B_n}{n!} \left(\frac{t}{2}\right)^n - \sum_{n=0}^{\infty} \frac{B_n}{n!} t^n = - \sum_{n=0}^{\infty} \frac{B_n}{n!} (1-2^{-n+1}) t^n,$$

$$B_n \left(\frac{1}{2}\right) = -(1-2^{-n+1}) B_n,$$

$$y_{2k} \left(\frac{1}{2}\right) = B_{2k} \left(\frac{1}{2}\right) - B_{2k} = -(2-2^{-2k+1}) B_{2k}.$$

Таким образом, внутри $[0, 1]$ $y_{2k}(x)$ сохраняет знак, противоположный знаку B_{2k} :

$$y_{2k}(x) B_{2k} < 0 \quad (0 < x < 1). \quad (\text{II.14})$$

В частности, так как B_{2k} и B_{2k+2} имеют противоположные знаки, $y_{2k}(x)$ и $y_{2k+2}(x)$ также будут иметь внутри $[0, 1]$ противоположные знаки.

§ 3. ПЕРИОДИЧЕСКИЕ ФУНКЦИИ, СВЯЗАННЫЕ С МНОГОЧЛЕНАМИ БЕРНУЛЛИ

Определим 1-периодические функции $B_n^*(x)$ равенствами

$$B_n^*(x) = B_n(x) \quad (0 \leq x < 1) \quad \text{и} \quad B_n^*(x+1) = B_n^*(x) \quad (-\infty < x < \infty),$$

$$B_0^*(x) \equiv 1.$$

Так как $B_1(x) = x - \frac{1}{2}$, то $B_1^*(x)$ есть разрывная функция, имеющая скачок -1 в целых точках. При $n > 1$ $B_n(1) = B_n(0)$ и $B_n^*(x)$ есть непрерывная периодическая функция.

Построим тригонометрические ряды для $B_n^*(x)$ на $[0, 1]$.

$$B_n^*(x) = \frac{1}{2} a_0^{(n)} + \sum_{m=1}^{\infty} [a_m^{(n)} \cos 2\pi m x + b_m^{(n)} \sin 2\pi m x], \quad (\text{II.15})$$

$$a_m^{(n)} = 2 \int_0^1 B_n^*(x) \cos 2\pi m x dx = 2 \int_0^1 B_n(x) \cos 2\pi m x dx,$$

$$b_m^{(n)} = 2 \int_0^1 B_n^*(x) \sin 2\pi m x dx = 2 \int_0^1 B_n(x) \sin 2\pi m x dx.$$

На основании известных в теории рядов Фурье теорем можно утверждать, что равенство (II.15) имеет место при всяких x , когда $n > 1$, ввиду непрерывности $B_n^*(x)$, и всюду, кроме целых точек, для $B_1^*(x)$, в целых же точках сумма ряда равна

$$-\frac{1}{2} [B_1^*(+0) + B_1^*(-0)] = -\frac{1}{2} \left[-\frac{1}{2} + \frac{1}{2} \right] = 0.$$

Подсчитаем коэффициенты $a_m^{(n)}$ и $b_m^{(n)}$. Пусть $n > 0$.

$$a_0^{(n)} = 2 \int_0^1 B_n(x) dx = \frac{2}{n+1} \int_0^1 B_{n+1}(x) dx = 0 \quad (n=1, 2, \dots).$$

Остановимся сначала на случае четного n ($n=2k$):

$$\begin{aligned} a_m^{(2k)} &= 2 \int_0^1 B_{2k}(x) \cos 2\pi m x dx = \\ &= 2 \left| \int_0^1 \frac{\sin 2\pi m x}{2\pi m} B_{2k}(x) - \frac{2k}{\pi m} \int_0^1 \sin 2\pi m x B_{2k-1}(x) dx. \right. \end{aligned}$$

Внеинтегральный член обращается в нуль. Повторное интегрирование по частям дает

$$a_m^{(2k)} = \frac{2k}{\pi m} \left| \int_0^1 \frac{\cos 2\pi m x}{2\pi m} B_{2k-1}(x) - \frac{k(2k-1)}{(\pi m)^2} \int_0^1 B_{2k-2}(x) \cos 2\pi m x dx. \right.$$

Если $k > 1$, внеинтегральный член обратится в нуль, так как

$$B_{2k-1}(1) = B_{2k-1}(0) = 0,$$

и получится

$$a_m^{(2k)} = -\frac{2k(2k-1)}{(2\pi m)^2} a_m^{(2k-2)}. \quad (\text{II.16})$$

Для $k=1$ интегральный член правой части равен нулю и

$$a_m^{(2)} = \frac{1}{(\pi m)^2}. \quad (\text{II.17})$$

Применение равенства (II.16) k раз даст совместно с (II.17) следующее значение для $a_m^{(2k)}$:

$$a_m^{(2k)} = (-1)^{k-1} \frac{2 \cdot (2k)!}{(2\pi m)^{2k}}.$$

Что же касается $b_m^{(2k)}$, то, так как

$$B_{2k}(1-x) = B_{2k}(x) \text{ и } \sin 2\pi m(1-x) = -\sin 2\pi mx,$$

для функции, стоящей под знаком интеграла в выражении $b_m^{(2k)}$, выполняется равенство

$$B_{2k}(1-x) \sin 2\pi m(1-x) = -B_{2k}(x) \sin 2\pi mx,$$

и поэтому $b_m^{(2k)} = 0$.

Ряд Фурье для $B_{2k}^*(x)$ будет

$$B_{2k}^*(x) = \frac{(-1)^{k-1} (2k)!}{2^{2k-1} \cdot \pi^{2k}} \sum_{m=1}^{\infty} \frac{\cos 2\pi mx}{m^{2k}}. \quad (\text{II.18})$$

При $x=0$ отсюда получаются равенства (II.6) для чисел Бернулли B_{2k} , приведенные нами в конце § 1.

Для $n=2k-1$ ($k=1, 2, \dots$) при помощи аналогичных вычислений найдем

$$B_{2k-1}^*(x) = \frac{(-1)^k (2k-1)!}{2^{2k-2} \pi^{2k-1}} \sum_{m=1}^{\infty} \frac{\sin 2\pi mx}{m^{2k-1}}.$$

§ 4. ПРЕДСТАВЛЕНИЕ ПРОИЗВОЛЬНОЙ ФУНКЦИИ ПРИ ПОМОЩИ МНОГОЧЛЕНОВ БЕРНУЛЛИ

Теорема 1. Если функция $f(x)$ имеет на $[0, 1]$ непрерывную производную порядка ν ($\nu \geq 1$), тогда при $0 \leq x \leq 1$ верно равенство

$$\begin{aligned} f(x) = & \int_0^1 f(t) dt + \sum_{k=1}^{\nu-1} \frac{B_k(x)}{k!} [f^{(k-1)}(1) - f^{(k-1)}(0)] - \\ & - \frac{1}{\nu!} \int_0^1 f^{(\nu)}(t) [B_{\nu}^*(x-t) - B_{\nu}^*(x)] dt. \end{aligned} \quad (\text{II.19})$$

Доказательство. Преобразуем интеграл

$$\rho_v(x) = \rho_v = \frac{1}{v!} \int_0^1 B_v^*(x-t) f^{(v)}(t) dt. \quad (\text{II.20})$$

Пусть $v > 1$. Выполним интегрирование по частям:

$$\rho_v = \left[\frac{B_v^*(x-t)}{v!} f^{(v-1)}(t) - \frac{1}{v!} \int_0^1 f^{(v-1)}(t) \frac{d}{dt} B_v^*(x-t) dt \right].$$

Так как

$$B_v^*(x-1) = B_v^*(x) = B_v(x) \quad \text{и} \quad \frac{d}{dt} B_v^*(x-t) = -v B_{v-1}^*(x-t),$$

то

$$\rho_v = -\frac{B_v(x)}{v!} [f^{(v-1)}(1) - f^{(v-1)}(0)] + \rho_{v-1}$$

и после $(v-1)$ -кратного применения этого соотношения получим

$$\rho_v = \sum_{k=2}^v \frac{B_k(x)}{k!} [f^{(k-1)}(1) - f^{(k-1)}(0)] + \rho_1.$$

Напомним, что $B_1^*(x)$ имеет разрывы со скачком -1 в целых точках и всюду между этими точками производная от $B_1^*(x)$ равна 1. Предположим, что $0 < x < 1$. Тогда

$$\begin{aligned} \rho_1(x) &= \int_0^x f'(t) B_1^*(x-t) dt + \int_x^1 f'(t) B_1^*(x-t) dt = B_1^*(+0) f(x) - \\ &- B_1^*(x) f(0) + \int_0^x f(t) dt + B_1^*(x-1) f'(1) - B_1^*(-0) f(x) + \int_x^1 f(t) dt = \\ &= [B_1^*(+0) + B_1^*(-0)] f(x) + B_1(x) [f(1) - f(0)] + \int_0^1 f(t) dt. \end{aligned}$$

Но

$$B_1^*(+0) = -\frac{1}{2}, \quad B_1^*(-0) = \frac{1}{2},$$

и, стало быть,

$$\rho_1 = -f(x) + B_1(x) [f(1) - f(0)] + \int_0^1 f(t) dt.$$

Для ρ_v ($v=1, 2, \dots$) получится представление, лишь формой записи отличающееся от (II.19):

$$\begin{aligned} & \frac{1}{v!} \int_0^1 f^{(v)}(t) B_v^*(x-t) dt = -f(x) + \\ & + \sum_{k=1}^v \frac{B_k(x)}{k!} [f^{(k-1)}(1) - f^{(k-1)}(0)] + \int_0^1 f(t) dt. \end{aligned}$$

Равенство (II.19) доказано нами в предположении, что x лежит внутри $[0, 1]$, но оно верно и для замкнутого отрезка $0 \leq x \leq 1$, так как все величины, в него входящие, являются непрерывными функциями x .

Когда функция $f(x)$ задана на произвольном конечном отрезке $[a, b]$ и имеет там непрерывную производную порядка v , ее разложение по многочленам Бернулли получится из (II.19) с помощью линейного преобразования аргумента. Введем новую переменную, положив $x = a + h\xi$ ($h = b - a$, $0 \leq \xi \leq 1$). К функции $\varphi(\xi) = f(a + h\xi)$ применим равенство (II.19):

$$\begin{aligned} \varphi(\xi) &= \int_0^1 \varphi(\tau) d\tau + \sum_{k=1}^{v-1} \frac{B_k(\xi)}{k!} [\varphi^{(k-1)}(1) - \varphi^{(k-1)}(0)] - \\ & - \frac{1}{v!} \int_0^1 \varphi^{(v)}(\tau) [B_v^*(\xi - \tau) - B_v^*(\xi)] d\tau \end{aligned}$$

и вернемся к первоначальному аргументу x и функции f , приняв во внимание соотношения

$$\varphi(\tau) = f(a + \tau h), \quad t = a + \tau h, \quad dt = h d\tau,$$

$$\varphi^{(k)}(\xi) = \frac{d^k}{d\xi^k} \varphi(\xi) = h^k \frac{d^k}{dx^k} f(a + h\xi) = h^k f^{(k)}(x).$$

Получим

$$\begin{aligned} f(x) &= \frac{1}{h} \int_a^b f(t) dt + \sum_{k=1}^{v-1} \frac{h^{k-1}}{k!} B_k\left(\frac{x-a}{h}\right) [f^{(k-1)}(b) - f^{(k-1)}(a)] - \\ & - \frac{h^{v-1}}{v!} \int_a^b f^{(v)}(t) \left[B_v^*\left(\frac{x-t}{h}\right) - B_v^*\left(\frac{x-a}{h}\right) \right] dt, \quad h = b - a. \quad (\text{II.21}) \end{aligned}$$

ДОБАВЛЕНИЕ III

АЛГЕБРАИЧЕСКИЕ МНОГОЧЛЕНЫ НАИЛУЧШЕГО ПРИБЛИЖЕНИЯ

Рассмотрим множество H_n алгебраических многочленов, степень которых не больше n :

$$P_n(x) = c_0x^n + c_1x^{n-1} + \dots + c_n.$$

Коэффициенты c_0, c_1, \dots, c_n есть любые действительные числа. В частности, c_0 может равняться нулю и степень многочлена P тогда будет меньше n . Очевидно, $H_0 \subset H_1 \subset H_2 \subset \dots$

Пусть на конечном отрезке $[a, b]$ дана непрерывная функция f . Возьмем многочлен P_n с определенными коэффициентами. Отклонение P_n от f на отрезке $[a, b]$ характеризуется величиной

$$\Delta(P_n) = \max_{a \leq x \leq b} |f(x) - P_n(x)|,$$

она зависит от выбора многочлена P_n и является функцией его коэффициентов c_0, c_1, \dots, c_n . Величина $\Delta(P_n)$ — неотрицательная и имеет неотрицательную точную нижнюю границу, когда P_n пробегает все множество H_n :

$$E_n = \inf_{P_n \in H_n} \Delta(P_n). \quad (\text{III.1})$$

Но можно показать, что E_n достигается и является минимальным значением $\Delta(P_n)$, так как существует такой многочлен P_n^* , для которого $\Delta(P_n^*) = E_n$. Мы не будем приводить арифметического доказательства существования P_n^* и ограничимся тем, что выясним наглядную сторону вопроса.

Величину E_n называют наименьшим отклонением многочленов из H_n от f или наилучшим приближением f многочленами степени n . $P_n^*(x)$ называется многочленом наилучшего приближения из H_n . Единственность такого многочлена будет установлена позже.

Отклонение $\Delta(P_n)$ есть, очевидно, непрерывная функция коэффициентов c_k . Введем $(n+1)$ -мерное числовое пространство R_{n+1} и совокупность значений коэффициентов c_k ($k=0, 1, \dots, n$) будем рассматривать как точку R_{n+1} . Так как c_k могут иметь произвольные значения, то областью определения $\Delta(P_n)$ будет все бесконечное пространство R_{n+1} . В этом состоит трудность задачи, так как нужно доказать существование минимума функции ΔP , непрерывной и заданной всюду в R_{n+1} . Но можно легко убедиться в том, что для нахождения минимума $\Delta(P_n)$ достаточно рассмотреть только ограниченную часть R_{n+1} .

При $n=0$, когда строится наилучшее приближение f постоянной величиной, задача решается просто. Если $m = \min_{[a, b]} f$ и $M = \max_{[a, b]} f$, то наимень-

шее отклонение от f имеет постоянная $P_0^* = \frac{M+m}{2}$. Здесь $\Delta(P_0^*) = \frac{M-m}{2}$.

Геометрически это вполне очевидно, так как если прямую $y = \frac{1}{2}(M+m)$ сдвинуть вверх или вниз, то ее наибольшее отклонение от линии $y=f(x)$ увеличится.

Множество многочленов H_n с увеличением n расширяется, поэтому при росте n точная нижняя граница отклонений не может увеличиваться, в частности, $E_n \leq E_0$ при $n > 0$.

График многочлена наилучшего приближения

$$P_n^*(x) = c_0^* x^n + c_1^* x^{n-1} + \dots + c_n^*,$$

если он существует, принадлежит замкнутой полосе между линиями $y=f(x)+E_n$, $y=f(x)-E_n$ и прямыми $x=a$, $x=b$.

Взяв произвольное $\alpha > 0$ и приняв во внимание неравенства

$$m \leq f(x) \leq M, \quad E_n \leq \frac{M-m}{2},$$

мы можем сказать, что значения $P_n^*(x)$, так же как и значения многочленов $P_n(x)$, близких к нему, т. е. имеющих коэффициенты c_k , достаточно близкие к c_k^* , лежат в границах

$$m - \frac{1}{2}(M-m) - \alpha \leq P_n(x) \leq M + \frac{1}{2}(M-m) + \alpha. \quad (\text{III.2})$$

Поэтому при изучении вопроса о существовании $P_n^*(x)$ мы можем не брать все многочлены, входящие в H_n , а ограничиться только теми, для которых выполняются неравенства (III.2). Нам осталось еще показать,

что коэффициентам c_k ($k=0, 1, \dots, n$) таких многочленов отвечает ограниченное множество точек в R_{n+1} . Возьмем на $[a, b]$ $n+1$ попарно разных узлов x_0, x_1, \dots, x_n и закрепим их. Интерполируем $P_n(x)$ по значениям во взятых узлах. Так как P_n имеет степень не выше n , интерполирование будет точным и

$$P_n(x) = \sum_{i=0}^n l_i(x) P_n(x_i), \quad l_i(x) = \frac{\omega(x)}{(x-x_i)\omega'(x_i)},$$

$$\omega(x) = (x-x_0)(x-x_1)\dots(x-x_n).$$

Если разложить многочлены влияния узлов $l_i(x)$ по степеням x :

$$l_i(x) = \sum_{j=0}^n d_{ij} x^j,$$

для коэффициентов c_k мы получим следующие выражения их через значения $P_n(x_k)$:

$$c_k = \sum_{i=0}^n d_{ik} P_n(x_i),$$

где d_{ik} зависят лишь от фиксированных узлов x_0, x_1, \dots, x_n и не зависят от значений $P_n(x_i)$. Теперь ясно, что если все значения многочлена $P_n(x)$ ограничены, то будут ограничены и все коэффициенты c_k . При решении вопроса о достижимости $\inf(\Delta P_n) = E_n$ мы можем считать, что ΔP_n , как функция коэффициентов c_0, c_1, \dots, c_n , задается в ограниченной замкнутой области R_{n+1} . Тогда ясно, что, ввиду непрерывной зависимости ΔP_n от c_k ($k=0, 1, \dots, n$), точная нижняя граница достигается и многочлен наилучшего приближения $P_n^*(x)$ существует.

Обратимся к теореме Чебышева о точках наибольшего отклонения P_n^* от f . Для сокращения обозначений условимся многочлен наилучшего приближения обозначать одной буквой P , отбрасывая знаки n и $*$. Для него

$$\max_x |P(x) - f(x)| = E_n$$

и существует хотя бы одна точка x_0 , в которой

$$|P(x_0) - f(x_0)| = E_n.$$

Такую точку называют точкой наибольшего отклонения или, кратко, (е)-точкой.

График многочлена $P(x)$ лежит между линиями $y=f(x)+E_n$ и $y=f(x)-E_n$. В точке x_0 он касается либо верхней линии, либо нижней.

Всякую точку x_0 , в которой он касается верхней линии, называют $(+)$ -точкой наибольшего отклонения или сокращенно $(+)$ -точкой, и аналогично всякую точку, где график многочлена касается нижней ограничивающей линии, называют $(-)$ -точкой.

Очевидно, должны существовать как $(+)$ -точки, так и $(-)$ -точки, так как если бы график многочлена P не касался, например, нижней линии $y=f(x)-E_n$, то мы могли бы сдвинуть график незначительно вниз, добавляя к многочлену $P(x)$ малую отрицательную постоянную, и получить многочлен, график которого лежит в более узкой полосе около линии $y=f(x)$, что противоречит тому, что $P(x)$ есть многочлен наименьшего отклонения от f .

Как оказывается, точек наибольшего отклонения на $[a, b]$ больше, чем две.

Теорема 1. На отрезке $[a, b]$ существует последовательность $n+2$ (e) -точек

$$x_1 < x_2 < \dots < x_{n+2},$$

которые попеременно есть $(+)$ -точки и $(-)$ -точки.

Доказательство. Такую последовательность (e) -точек часто называют чебышевским альтернансом.

Доказательство теоремы основано на простой мысли: если чебышевский альтернанс из $n+2$ точек отсутствует и можно построить альтернанс самое большее из m точек при $m < n+2$, то отклонение P от f можно уменьшить, вычитая из P многочлен ρ , имеющий степень не выше n и надлежащим образом подобранный. Для этого нужно, чтобы многочлен ρ имел во всех (e) -точках тот же знак, что и разность $P-f$. Если эта разность меняет знак меньше чем $n+1$ раз, то многочленом степени n такому требованию можно удовлетворить. При этом, если ρ умножить на достаточно малый положительный множитель λ , новый многочлен $P(x) - \lambda\rho(x)$ во всех других точках будет отклоняться от f меньше, чем на E_n .

Это показало бы, что при отсутствии $(n+2)$ -членного альтернанса многочлен P не может иметь наименьшего отклонения от f .

Для построения ρ недостаточно указанных наглядных соображений, и оно потребует некоторых численных расчетов.

Отрезок $[a, b]$ разделим точками

$$a = u_0 < u_1 < \dots < u_s = b$$

на столь малые части, чтобы в каждой из них изменение функции *) $P-f$ было меньше $\frac{1}{2} E_n$. Каждую часть $u_k \leq x \leq u_{k+1}$, содержащую хоть одну (e) -точку, будем называть (e) -сегментом. На каждом таком сегменте

*) Под изменением функции $\varphi(x)$ на отрезке $[c, d]$ понимают $\sup_{c \leq x, x' \leq d} |\varphi(x) - \varphi(x')|$.

разность $P-f$ не обращается в нуль и сохраняет знак. Поэтому на (е)-сегменте могут лежать либо только $(+)$ -точки, тогда мы будем его называть $(+)$ -сегментом, либо только $(-)$ -точки — и мы будем называть его $(-)$ -сегментом. Заметим, что на каждом $(+)$ -сегменте разность $P-f$ положительна и на каждом $(-)$ -сегменте она отрицательна.

Разделим теперь (е)-сегменты на группы, позволяющие подсчитать наименьшее возможное число перемен знаков у разности $P-f$.

Перенумеруем (е)-сегменты слева направо:

$$d_1, d_2, \dots, d_N.$$

Разобьем их на группы по приведенной ниже схеме. Для определенности записи схемы мы считаем, что d_1 есть $(+)$ -сегмент.

$$\left. \begin{array}{ll} d_1, d_2, \dots, d_{k_1} & [(+)\text{-сегменты}], \\ d_{k_1+1}, d_{k_1+2}, \dots, d_{k_2} & [(-)\text{-сегменты}], \\ \dots & \dots \\ d_{k_{m-1}+1}, \dots, d_{k_m} & [(-1)^{m-1}\text{-сегменты}]. \end{array} \right\} \quad (\text{III.3})$$

В схеме показаны m групп. Каждая из них содержит по меньшей мере один (е)-сегмент. Для доказательства теоремы достаточно установить неравенство $m \geq n+2$.

Допустим противоположное: $m < n+2$ и убедимся, что такое допущение приводит к противоречию с тем, что P есть многочлен наименьшего отклонения от f . На сегментах d_{k_i} и d_{k_i+1} ($i=1, 2, \dots, m-1$) разность $P-f$ имеет противоположные знаки, такие сегменты не могут иметь общие концы и должны быть разделены между собой не-(е)-сегментами. Можно выбрать точку z_i ($i=1, \dots, m-1$), лежащую справа от d_{k_i} и слева от d_{k_i+1} .

Построим многочлен

$$\rho(x) = (z_1 - x)(z_2 - x) \dots (z_{m-1} - x).$$

Он имеет степень $m-1 \leq n$. На сегментах первой группы (III.3) ρ положителен, как и разность $P-f$, на сегментах второй группы ρ имеет отрицательные значения, подобно $P-f$ и т. д. На всех (е)-сегментах ρ и $P-f$ имеют одинаковые знаки.

На всех не-(е)-сегментах величина $|P-f|$ имеет значение, строго меньшее E_n . Пусть там выполняется неравенство $|P-f| \leq E' < E_n$.

Обозначим $\max_{[a, b]} |\rho(x)| = R$ и выберем положительное число λ настолько малым, чтобы было

$$\lambda R < E_n - E' \quad \text{и} \quad \lambda R < \frac{1}{2} E_n. \quad (\text{III.4})$$

Наконец, рассмотрим многочлен

$$Q(x) = P(x) - \lambda p(x),$$

имеющий степень не больше n , и покажем, что он отклоняется от f меньше, чем на E_n .

В самом деле, на всяком не-(е)-сегменте

$$|Q-f| \leq |P-f| + \lambda |\rho| \leq E' + \lambda R < E' + (E_n - E') = E_n.$$

Если же точка x лежит на (е)-сегменте, то, ввиду того что там $P-f$ и ρ имеют одинаковые знаки и $|P-f| > \frac{1}{2} E_n$, $|\lambda \rho| < \frac{1}{2} E_n$ и, кроме того, $\rho(x) \neq 0$, будет:

$$|Q-f| = |P - \lambda \rho - f| = |P-f| - \lambda |\rho| \leq E_n - \lambda |\rho| < E_n$$

и P , следовательно, не является многочленом наилучшего приближения.

Докажем теперь единственность многочлена наилучшего приближения.

Теорема 2. Среди многочленов степени не выше n существует единственный, имеющий наименьшее отклонение от f .

Доказательство. Предположим, что существуют в H_n два многочлена P и Q , которые имеют наименьшее отклонение от f . Для них выполняются неравенства

$$-E_n \leq P-f \leq E_n,$$

$$-E_n \leq Q-f \leq E_n.$$

Если сложить их почленно и результат разделить на 2, получится новое неравенство

$$-E_n \leq \frac{P+Q}{2} - f \leq E_n,$$

говорящее, что $R = \frac{1}{2} (P+Q)$ также есть многочлен наилучшего приближения. По теореме 1 для R должны существовать $n+2$ точки, где наибольшее отклонение достигается. Пусть это будут точки x_1, x_2, \dots, x_{n+2} .

Если x_k есть одна из (+)-точек, то

$$R(x_k) - f(x_k) = \frac{P(x_k) - f(x_k)}{2} + \frac{Q(x_k) - f(x_k)}{2} = E_n.$$

Но $Q(x_k) - f(x_k) \leq E_n$ и, следовательно,

$$\frac{P(x_k) - f(x_k)}{2} + \frac{1}{2} E_n \geq E_n.$$

Значит,

$$P(x_k) - f(x_k) \geq E_n.$$

Но здесь возможен только знак равенства:

$$P(x_k) - f(x_k) = E_n.$$

Аналогично

$$Q(x_k) - f(x_k) = E_n$$

и, стало быть,

$$Q(x_k) = P(x_k).$$

Сходным путем доказывается, что P и Q имеют одинаковые значения и в $(-)$ -точках. Таким образом, два многочлена степени не выше n оказались совпадающими в $n+2$ точках. Последнее же может быть только в случае их тождественного равенства.

Добавление IV

НЕКОТОРЫЕ СВЕДЕНИЯ ОБ УРАВНЕНИЯХ В КОНЕЧНЫХ РАЗНОСТЯХ

§ 1. УРАВНЕНИЯ В КОНЕЧНЫХ РАЗНОСТЯХ ПРОИЗВОЛЬНОГО ВИДА

Пусть функция $y(x)$ задана в некоторой области. Для определенности предположим, что область ее задания есть полуось $0 \leq x < \infty$. Возьмем сетку равноотстоящих точек $x + kh$ с шагом $h > 0$ и рассмотрим конечные разности функции $y(x)$:

$$\Delta y(x) = y(x+h) - y(x), \dots, \Delta^p y(x) = \Delta^{p-1} y(x+h) - \Delta^{p-1} y(x).$$

Уравнение вида

$$\Phi(x, y(x), \Delta y(x), \dots, \Delta^p y(x)) = 0 \quad (\text{IV.1})$$

называется уравнением в конечных разностях порядка p . Под $y(x)$ здесь понимается функция, подлежащая нахождению. $\Phi(x, y_0, y_1, \dots, y_p)$ есть заданная функция, определенная в некоторой области изменения своих аргументов $(x, y_0, y_1, \dots, y_p)$.

При помощи известных выражений конечных разностей через значения функции

$$\Delta y(x) = y(x+h) - y(x), \Delta^2 y(x) = y(x+2h) - 2y(x+h) + y(x), \dots,$$

$$\Delta^p y(x) = y(x+ph) - \frac{p}{1} y(x+(p-1)h) + \dots + (-1)^p y(x)$$

уравнению (IV.1) можно придать форму

$$\Psi(x, y(x), y(x+h), \dots, y(x+ph)) = 0. \quad (\text{IV.2})$$

Если же воспользоваться соотношениями

$$y(x+h) = y(x) + \Delta y(x), \quad y(x+2h) = y(x) + 2\Delta y(x) + \Delta^2 y(x), \dots,$$

$$y(x+ph) = y(x) + \frac{p}{1} \Delta y(x) + \dots + \Delta^p y(x),$$

уравнение (IV.2) можно легко привести к виду (IV.1).

Ниже мы будем рассматривать уравнение в конечных разностях в форме (IV.2).

В нем независимая переменная x может иметь, вообще говоря, любые неотрицательные значения. Для нас же достаточно считать, что x принимает значения вида $x = nh$ ($n = 0, 1, 2, \dots$). Обозначая $y(kh) = y_k$, мы будем записывать разностное уравнение в виде

$$F(n, y_n, y_{n+1}, \dots, y_{n+p}) = 0 \quad (n = 0, 1, \dots). \quad (\text{IV.3})$$

Условимся говорить, что уравнение имеет нормальную форму, если оно решено относительно значения y с наибольшим индексом:

$$y_{n+p} = f(n, y_n, y_{n+1}, \dots, y_{n+p-1}). \quad (\text{IV.4})$$

Его можно рассматривать как рекурсионное соотношение весьма частного вида, дающего явное выражение значения y_{n+p} функции y через n и p предшествующих значений y_{n+p-1}, \dots, y_n .

Если считать y_0, y_1, \dots, y_{p-1} известными, то, полагая $n=0$, при помощи (IV.4) мы вычислим $y_p = f(0, y_0, y_1, \dots, y_{p-1})$. Полагая $n=1$, найдем $y_{p+1} = f(1, y_1, y_2, \dots, y_p)$ и т. д. Вычисления можно продолжать либо неограниченно далеко, либо до тех пор, когда точка $(n, y_n, y_{n+1}, \dots, y_{n+p-1})$ выйдет из области определения функции f .

Что же касается начальных значений y_0, y_1, \dots, y_{p-1} , то они остаются произвольными и им можно придавать любые значения также из области определения f .

§ 2. ЛИНЕЙНЫЕ УРАВНЕНИЯ

После кратких пояснений предыдущего параграфа остановимся более подробно на линейных разностных уравнениях. Это уравнения, линейные относительно значений неизвестной функции y_n :

$$L(y_n) = a_0(n)y_{n+p} + a_1(n)y_{n+p-1} + \dots + a_p(n)y_n = f(n). \quad (\text{IV.5})$$

Коэффициенты $a_k(n)$ и свободный член $f(n)$ могут быть произвольными функциями целочисленного аргумента n . Мы будем для простоты считать, что они определены для всех неотрицательных значений n . Если $a_0(n) \neq 0$ для некоторого значения n , то из уравнения (IV.5) можно найти y_{n+p} в форме линейной функции от значения $f(n)$ свободного члена и p предшествующих значений функции y : y_{n+p-1}, \dots, y_n . Поэтому, если $a_0(n)$ не обращается в нуль ни при каких значениях n (или на некотором отрезке $0 \leq n \leq N$), при помощи уравнения мы можем, решая его последовательно относительно y_p, y_{p+1}, \dots , найти значение y_n любого номера n ($n \geq p$) в виде линейной функции начальных значений y_0, y_1, \dots, y_{p-1} и значений $f(0), f(1), \dots, f(n-p)$ свободного члена:

$$y_n = \sum_{i=0}^{p-1} \Gamma_n^i y_i + \sum_{j=0}^{n-p} G_n^j f(j) \quad (n=p, p+1, \dots). \quad (\text{IV.6})$$

К рассмотрению такого представления решения вернемся несколькими страницами позже. Сейчас же нам необходимо отметить, что если $a_0(n) \neq 0$, то при всяких y_0, y_1, \dots, y_{p-1} уравнение имеет решение с такими начальными значениями и это решение единственное.

Уравнение (IV.5) называется однородным, если его свободный член $f(n)$ тождественно равен нулю. Такое уравнение имеет вид

$$L(z_n) = a_0(n)z_{n+p} + a_1(n)z_{n+p-1} + \dots + a_p(n)z_n = 0. \quad (\text{IV.7})$$

Для его решений верно, очевидно, следующее утверждение.

Если $z_n^{(1)}, z_n^{(2)}, \dots, z_n^{(k)}$ есть решения однородного уравнения (IV.7), то их линейная комбинация с произвольными постоянными коэффициентами C_j ($j=1, \dots, k$)

$$z_n = C_1 z_n^{(1)} + C_2 z_n^{(2)} + \dots + C_k z_n^{(k)}$$

есть также решение однородного уравнения (IV.7).

Покажем сейчас, что для построения всякого решения однородного уравнения (IV.7) достаточно знать p его частных решений, обладающих свойством, о котором мы будем говорить.

Рассмотрим p решений однородного уравнения $z_n^{(i)}$ ($i=1, 2, \dots, p$). Говорят, что эти решения образуют фундаментальную систему, если определитель, составленный из их начальных значений, отличен от нуля:

$$W_p = \begin{vmatrix} z_0^{(1)} & z_0^{(2)} & \dots & z_0^{(p)} \\ z_1^{(1)} & z_1^{(2)} & \dots & z_1^{(p)} \\ \dots & \dots & \dots & \dots \\ z_{p-1}^{(1)} & z_{p-1}^{(2)} & \dots & z_{p-1}^{(p)} \end{vmatrix} \neq 0. \quad (\text{IV.8})$$

Такое название связано с тем обстоятельством, что всякое решение однородного уравнения есть линейная комбинация решений, образующих фундаментальную систему. Действительно, каждое решение определяется начальными значениями. Обозначим их a_0, a_1, \dots, a_{p-1} .

Образует линейную комбинацию

$$z_n = C_1 z_n^{(1)} + C_2 z_n^{(2)} + \dots + C_p z_n^{(p)}. \quad (\text{IV.9})$$

При произвольных C_k она является решением однородного уравнения.

Нам осталось установить, что C_k можно избрать так, чтобы решение z_n имело заданные заранее начальные значения. Полагая $n=0, 1, \dots, p-1$ и приравнявая соответствующие им значения z_n заданным числам a_k , получим для C_k линейную систему:

$$\begin{aligned} z_0 &= C_1 z_0^{(1)} + C_2 z_0^{(2)} + \dots + C_p z_0^{(p)} = a_0, \\ z_1 &= C_1 z_1^{(1)} + C_2 z_1^{(2)} + \dots + C_p z_1^{(p)} = a_1, \\ &\vdots \\ z_{p-1} &= C_1 z_{p-1}^{(1)} + C_2 z_{p-1}^{(2)} + \dots + C_p z_{p-1}^{(p)} = a_{p-1}. \end{aligned}$$

Определитель ее совпадает с W_p и, так как $W_p \neq 0$, из системы всегда могут быть найдены и при этом единственным образом коэффициенты C_1, \dots, C_p .

Линейную комбинацию (IV.9) часто называют общим решением однородного разностного уравнения (IV.7).

Понятие фундаментальной системы тесно связано с понятием линейной независимости решений однородного уравнения. Будем вновь считать, что нам известно k решений $z_n^{(1)}, \dots, z_n^{(k)}$. Эти решения называются линейно зависимыми, если существуют такие постоянные величины C_1, \dots, C_k , не все равные нулю, что при всяких n выполняется равенство

$$C_1 z_n^{(1)} + C_2 z_n^{(2)} + \dots + C_k z_n^{(k)} = 0.$$

Если же такое равенство при всех n может выполняться только в том случае, когда все C_i ($i=1, \dots, k$) равны нулю, решения называются линейно независимыми.

Можно просто проверить, что среди решений однородного уравнения (IV.5) существует не более чем p линейно независимых. В самом деле, пусть $z_n^{(1)}, \dots, z_n^{(p)}, z_n^{(p+1)}$ есть произвольные $p+1$ таких решений, образуем из них линейную комбинацию с постоянными коэффициентами

$$z_n' = C_1 z_n^{(1)} + \dots + C_p z_n^{(p)} + C_{p+1} z_n^{(p+1)}.$$

При любых C_i это есть решение однородного уравнения. Покажем, что C_i можно выбрать так, чтобы это решение имело нулевые начальные значения, при этом не все C_i будут равны нулю. Для этого нужно выполнить систему p уравнений

$$C_1 z_n^{(1)} + \dots + C_p z_n^{(p)} + C_{p+1} z_n^{(p+1)} = 0 \quad (n=0, 1, \dots, p-1).$$

Это есть однородная система p уравнений с $p+1$ неизвестными C_k ($k=1, \dots, p+1$). Такая система всегда имеет ненулевое решение.

Возьмем любое из таких решений. Соответствующая ему линейная комбинация z_n' , как решение уравнения (IV.7) с начальными значениями, равными нулю, будет равна нулю при всяких n , а отсюда следует,

что решения $z_n^{(1)}, \dots, z_n^{(p+1)}$ являются линейно зависимыми.

Возьмем теперь p решений $z_n^{(1)}, \dots, z_n^{(p)}$ уравнения (IV.7) и покажем, что для линейной зависимости их необходимо и достаточно, чтобы определитель W_p (IV.6), составленный из их начальных значений, был равен нулю.

Образует из них линейную комбинацию

$$z_n = C_1 z_n^{(1)} + \dots + C_p z_n^{(p)}.$$

Если решения линейно зависимы, существуют C_k ($k=1, \dots, p$), не все равные нулю и такие, что $z_n \equiv 0$. В частности, $z_n = 0$ при $n=0, 1, \dots, p-1$. Это даст для C_k однородную систему

$$C_1 z_n^{(1)} + \dots + C_p z_n^{(p)} = 0 \quad (n=0, 1, \dots, p-1)$$

с определителем W_p . Так как системе удовлетворяют C_k , среди которых есть не равные нулю, определитель системы должен равняться нулю: $W_p = 0$.

Наоборот, если $W_p = 0$, то последняя однородная система имеет ненулевое решение. Соответствующая ему комбинация z_n будет иметь нулевые начальные значения и будет, следовательно, равна нулю при всех n , что говорит о линейной зависимости взятых решений.

Последние рассуждения показывают, что следующие утверждения являются равносильными:

1) решения $z_n^{(1)}, \dots, z_n^{(p)}$ уравнения (IV.7) образуют фундаментальную систему и

2) эти решения линейно независимы.

Возвратимся к произвольному линейному уравнению (IV.5). Легко проверить, что разность между двумя решениями y_n и $y_n^{(0)}$ неоднородного уравнения (IV.5) есть решение однородного уравнения (IV.7):

$$L(y_n - y_n^{(0)}) = L(y_n) - L(y_n^{(0)}) = f(n) - f(n) = 0.$$

Если эту разность обозначить z_n , то $y_n = y_n^{(0)} + z_n$. Верно также утверждение: если $y_n^{(0)}$ есть решение неоднородного уравнения (IV.5) и z_n есть

решение однородного уравнения (IV.7), то $y_n = y_n^{(0)} + z_n$ есть решение неоднородного уравнения (IV.5). В самом деле,

$$L(y_n) = L(y_n^{(0)}) + L(z_n) = f(n) + 0 = f(n).$$

Иными словами говоря, всякое решение y_n неоднородного уравнения (IV.5) представимо в виде

$$y_n = y_n^{(0)} + z_n = y_n^{(0)} + C_1 z_n^{(1)} + \dots + C_p z_n^{(p)}, \quad (\text{IV.10})$$

где $y_n^{(0)}$ есть некоторое решение неоднородного уравнения и z_n — решение однородного уравнения. Верно также, что если $y_n^{(0)}$ есть некоторое решение (IV.5) и z — любое решение (IV.7), то (IV.10) есть решение неоднородного уравнения (IV.5).

(IV.10) называют, ввиду изложенного, общим решением неоднородного уравнения.

Выше было указано выражение (IV.6) для произвольного решения y_n уравнения (IV.5) через начальные значения y_0, y_1, \dots, y_{p-1} и значения свободного члена $f(n)$. Это также есть одно из представлений общего решения, записанное лишь в иной форме, чем (IV.10), со специальным выбором фундаментальной системы

$$z_n^{(1)} = \Gamma_n^0, z_n^{(2)} = \Gamma_n^1, \dots, z_n^{(p)} = \Gamma_n^{p-1}$$

и частного решения неоднородного уравнения

$$y_n^{(0)} = \sum_{j=0}^{n-p} G_n^j f(j).$$

Выясним наглядный смысл коэффициентов Γ_n^i и G_n^j . Начнем с Γ_n^0 . Во-первых, $f(n) \equiv 0$ и будем, следовательно, уравнение считать однородным. В правой части (IV.6) при этом исчезнет вторая сумма.

Во-вторых, начальные значения выберем следующими:

$$y_0 = 1, y_1 = \dots = y_{p-1} = 0.$$

Тогда мы получим $y_n = \Gamma_n^0$. Это позволяет сказать, что Γ_n^0 есть решение однородного уравнения (IV.7), удовлетворяющее начальным условиям $\Gamma_0^0 = 1, \Gamma_1^0 = 0, \dots, \Gamma_{p-1}^0 = 0$. Такое решение учитывает влияние, которое оказывает на y_n начальное значение y_0 . Аналогично Γ_n^1 будет решением однородного уравнения с начальными значениями

$$\Gamma_0^i = 0, \Gamma_1^i = 1, \Gamma_2^i = \dots = \Gamma_{p-1}^i = 0.$$

Оно будет учитывать влияние на y_n начального значения y_2 . Сходный смысл имеют прочие Γ_n^i .

Рассмотренные коэффициенты Γ_n^i называют функциями Грина или функциями влияния начальных значений.

Перейдем теперь к выяснению роли G_n^i . С этой целью положим начальные значения y_i равными нулю. При этом исчезнет первая сумма справа в (IV.6). Затем фиксируем какое-либо значение индекса $n=i$ и положим $f(n)=0$ при $n \neq i$ и $f(i)=1$. Если воспользоваться символом Кронекера, такой выбор $f(n)$ можно записать в виде $f(n)=\delta_n^i$.

Считая $n \geq p+i$, мы получим для y_n равенство $y_n = G_n^i$, что дает возможность сказать, что G_n^i есть решение неоднородного уравнения

$$L(G_n^i) = \delta_n^i \quad (n = p, p+1, \dots) \quad (\text{IV.11})$$

с единственным, отличным от нуля значением свободного члена $f(i) = \delta_i^i = 1$. Начальные значения этого решения все равны нулю:

$$G_n^i = 0 \quad (n = 0, 1, \dots, p-1).$$

Коэффициент G_n^i называется функцией влияния или гриновой функцией значения $f(i)$ свободного члена.

§ 3. ЛИНЕЙНЫЕ УРАВНЕНИЯ С ПОСТОЯННЫМИ КОЭФФИЦИЕНТАМИ

Рассмотрим теперь линейные однородные уравнения, коэффициенты которых не зависят от n и являются постоянными:

$$L(z_n) = a_0 z_{n+p} + a_1 z_{n+p-1} + \dots + a_n z_n = 0 \quad (n = 0, 1, \dots). \quad (\text{IV.12})$$

Как мы покажем сейчас, построение общего решения здесь приводится к нахождению корней алгебраического уравнения степени p и к определению их кратностей.

Сделаем замену функции в уравнении, положив $z_n = \lambda^n u_n$, где λ есть некоторая постоянная величина, выбор которой будет сделан ниже.

Воспользовавшись выражением значений функции через конечные разности, получим

$$\begin{aligned}
z_n &= \lambda^n u_n, \\
z_{n+1} &= \lambda^{n+1} u_{n+1} = \lambda^{n+1} (u_n + \Delta u_n), \\
z_{n+2} &= \lambda^{n+2} u_{n+2} = \lambda^{n+2} (u_n + 2\Delta u_n + \Delta^2 u_n), \\
&\dots \\
z_{n+p} &= \lambda^{n+p} u_{n+p} = \lambda^{n+p} \left(u_n + \frac{p}{1!} \Delta u_n + \frac{p(p-1)}{2!} \Delta^2 u_n + \dots + \Delta^p u_n \right).
\end{aligned}$$

Подстановка в (IV.12) даст:

$$\begin{aligned}
L(\lambda^n u_n) &= \lambda^{n+p} a_0 \left[u_n + \frac{p}{1!} \Delta u_n + \frac{p(p-1)}{2!} \Delta^2 u_n + \dots \right] + \\
&+ \lambda^{n+p-1} a_1 \left[u_n + \frac{p-1}{1!} \Delta u_n + \dots \right] + \dots + \lambda^n a_p u_n = \lambda^n \varphi(\lambda) u_n + \\
&+ \lambda^{n+1} \varphi'(\lambda) \frac{\Delta u_n}{1!} + \lambda^{n+2} \varphi''(\lambda) \frac{\Delta^2 u_n}{2!} + \dots + \lambda^{n+p} \varphi^{(p)}(\lambda) \frac{\Delta^p u_n}{p!} = 0.
\end{aligned}$$

Здесь

$$\varphi(\lambda) = a_0 \lambda^p + a_1 \lambda^{p-1} + a_2 \lambda^{p-2} + \dots + a_p.$$

Предположим теперь, что λ есть корень уравнения $\varphi(\lambda) = 0$ и кратность его равна k . В этом случае

$$\varphi(\lambda) = \varphi'(\lambda) = \dots = \varphi^{(k-1)}(\lambda) = 0$$

и

$$L(\lambda^n u_n) = \lambda^{n+k} \varphi^{(k)}(\lambda) \frac{\Delta^k u_n}{k!} + \dots + \lambda^{n+p} \varphi^{(p)}(\lambda) \frac{\Delta^p u_n}{p!}.$$

Если, кроме того, u_n есть многочлен от n , степень которого не больше $k-1$, то будет $\Delta^k u_n = 0, \dots, \Delta^p u_n = 0$ и $L(\lambda^n u_n) = 0$.

Все это дает возможность высказать приводимое ниже заключение. Пусть уравнение

$$\varphi(\lambda) = a_0 \lambda^p + a_1 \lambda^{p-1} + \dots + a_p = 0 \quad (\text{IV.13})$$

имеет m различных корней $\lambda_1, \lambda_2, \dots, \lambda_m$, кратности которых равны соответственно k_1, k_2, \dots, k_m . Этим корням отвечают следующие частные решения уравнения:

$$\left. \begin{aligned}
\text{корню } \lambda_1: & \quad \lambda_1^n, n\lambda_1^n, n^2\lambda_1^n, \dots, n^{k_1-1}\lambda_1^n, \\
\text{корню } \lambda_2: & \quad \lambda_2^n, n\lambda_2^n, n^2\lambda_2^n, \dots, n^{k_2-1}\lambda_2^n, \\
& \dots \\
\text{корню } \lambda_m: & \quad \lambda_m^n, n\lambda_m^n, n^2\lambda_m^n, \dots, n^{k_m-1}\lambda_m^n.
\end{aligned} \right\} \quad (\text{IV.14})$$

Так как $k_1 + k_2 + \dots + k_m = p$, приведенных решений будет ровно p .

Почти очевидной является их линейная независимость. Чтобы проверить ее, достаточно составить из их начальных значений при $n = 0, 1, \dots, p-1$ определитель W_p и проверить неравенство его нулю. Связанные с этим вычисления не имеют принципиальных трудностей, но громоздки по выполнению. Поэтому мы не будем останавливаться на доказательстве того, что $W_p \neq 0$, и ограничимся только тем, что отметим некоторые частные случаи.

1. Если все корни уравнения $\varphi(\lambda) = 0$ являются однократными ($k_1 = k_2 = \dots = 1$), решениями однородного уравнения будут

$$\lambda_1^n, \lambda_2^n, \dots, \lambda_p^n. \quad (\text{IV.15})$$

Определитель W_p для них есть

$$W_p = \begin{vmatrix} 1 & 1 & \dots & 1 \\ \lambda_1 & \lambda_2 & \dots & \lambda_p \\ \cdot & \cdot & \cdot & \cdot \\ \lambda_1^{p-1} & \lambda_2^{p-1} & \dots & \lambda_p^{p-1} \end{vmatrix}.$$

Он является определителем Вандермонда и отличен от нуля. Поэтому решения (IV.14) линейно независимы и образуют фундаментальную систему.

2. Рассмотрим уравнение второго порядка

$$L(z_n) = a_0 z_{n+2} + a_1 z_{n+1} + a_2 z_n = 0.$$

Алгебраическое уравнение (IV.13) здесь будет квадратным:

$$\varphi(\lambda) = a_0 \lambda^2 + a_1 \lambda + a_2 = 0.$$

Пусть корни его есть λ_1 и λ_2 . Если $\lambda_1 \neq \lambda_2$, то решениями будут λ_1^n и λ_2^n .

Если же $\lambda_1 = \lambda_2$, то решения есть λ_1^n и $n\lambda_1^n$. Они, очевидно, линейно независимы.

3. Возвратимся к общему случаю (IV.14). Когда модули корней различны между собой ($|\lambda_i| \neq |\lambda_j|$, $i \neq j$), решения также будут очевидным образом независимы, так как они все будут иметь различные порядки роста при $n \rightarrow \infty$.

Предисловие	3
ГЛАВА 1. РЕШЕНИЕ ЧИСЛЕННЫХ УРАВНЕНИЙ	7
§ 1.1. О содержании задачи решения уравнений	—
§ 1.2. Метод итерации. Случай одного численного уравнения	10
§ 1.3. О задаче улучшения метода итерации. Некоторые видоизменения итерационного процесса	18
§ 1.4. Улучшение итерационного процесса при помощи преобразования заданного уравнения	28
§ 1.5. Понятие об общей теории метода итерации. Теорема о сжатых отображениях	34
§ 1.6. Метод итерации для систем уравнений	37
§ 1.7. Метод Ньютона. Случай одного численного уравнения	44
§ 1.8. Об уточнениях и изменениях метода Ньютона	56
§ 1.9. Операторные уравнения и метод Ньютона	73
§ 1.10. Метод Ньютона для систем уравнений	79
§ 1.11. Метод решения, основанный на возведении корней в степень	88
§ 1.12. Нахождение корней многочленов при помощи выделения множителей	96
Л и т е р а т у р а	102
ГЛАВА 2. РЕШЕНИЕ СИСТЕМ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ	103
§ 2.1. Некоторые сведения из линейной алгебры	104
2.1.1. Сходимость последовательностей векторов и матриц	—
2.1.2. Нормы векторов и матриц	—
2.1.3. Сходимость матричной геометрической прогрессии	115
§ 2.2. Итерационные методы	119
2.2.1. Основные разновидности итерационных процессов	120
2.2.2. Метод простой итерации	122
2.2.3. Метод Рундсона	132
2.2.4. Метод Зейделя и метод релаксации	136
§ 2.3. Методы исключения	150
2.3.1. Метод Гаусса	151
2.3.2. Метод оптимального исключения	155
2.3.3. Метод окаймления	158
2.3.4. Вычисление определителей	162
2.3.5. Обращение матриц	164
§ 2.4. Методы, основанные на разложениях матрицы	173
2.4.1. Метод квадратного корня	174
2.4.2. Метод отражений	181
2.4.3. Вычисление определителей	186
2.4.4. Обращение матриц	188
§ 2.5. Методы, основанные на построении вспомогательной системы векторов, ортогональных в некоторой метрике	189
2.5.1. Метод ортогонализации	—
2.5.2. Алгоритм Уилкинсона	197
2.5.3. Метод сопряженных градиентов	199
2.5.4. Вариант метода сопряженных градиентов	205
2.5.5. Метод скорейшего спуска	208
§ 2.6. Способы оценки погрешности приближенного решения системы	213
2.6.1. Обусловленность систем уравнений и матриц	214
2.6.2. Оценка погрешности	215
Л и т е р а т у р а	218

ГЛАВА 3. ВЫЧИСЛЕНИЕ СОБСТВЕННЫХ ЗНАЧЕНИЙ И СОБСТВЕННЫХ ВЕКТОРОВ МАТРИЦ	219
§ 3.1. О содержании задачи	—
§ 3.2. Метод А. Н. Крылова	224
3.2.1. Некоторые сведения из алгебры	225
3.2.2. Нахождение собственных значений матрицы	227
3.2.3. Вычисление собственных векторов матрицы	230
§ 3.3. Метод А. М. Данилевского	233
3.3.1. Построение собственного многочлена матрицы	234
3.3.2. Вычисление собственных векторов матрицы	239
§ 3.4. Другие методы получения собственного многочлена матрицы	241
3.4.1. Интерполяционный метод	—
3.4.2. Метод Леверье	243
3.4.3. Метод Д. К. Фаддеева	244
3.4.4. Метод окаймления	246
3.4.5. Эскалаторный метод	248
3.4.6. Метод ортогонализации	250
3.4.7. Метод Хессенберга	260
3.4.8. Метод Самуэльсона	265
§ 3.5. Итерационные методы нахождения собственных значений и собственных векторов матрицы	267
3.5.1. Степенной метод для вычисления наибольшего по модулю собственного значения матрицы и соответствующего собственного вектора	268
3.5.2. Вычисление всех собственных значений положительно определенной симметрической матрицы	279
3.5.3. Видоизменения степенного метода	282
3.5.4. Метод λ -разности	287
§ 3.6. Метод вращений	289
3.6.1. Случай вещественных симметрических матриц	292
3.6.2. Сходимость метода вращений	295
3.6.3. Случай эрмитовых матриц	301
§ 3.7. Уточнение собственных значений и принадлежащих им собственных векторов матриц и ускорение сходимости метода итерации при решении систем линейных алгебраических уравнений	304
3.7.1. Уточнение полной проблемы собственных значений	—
3.7.2. Уточнение отдельного собственного значения и принадлежащего ему собственного вектора	309
3.7.3. δ^2 -Процесс Эйткена	314
3.7.4. Метод М. К. Гавурина	317
3.7.5. Метод Л. А. Люстерника	320
Л и т е р а т у р а	324
ГЛАВА 4. ИНТЕРПОЛИРОВАНИЕ	325
§ 4.1. О содержании задачи интерполирования	—
4.1.1. Об интерполяционных приближениях	—
4.1.2. Остаток интерполирования	332
§ 4.2. Конечные разности и разностные отношения	336
4.2.1. Конечные разности	—
4.2.2. Разностные отношения, их свойства и связь с конечными разностями	338
§ 4.3. Алгебраическое интерполирование по значениям функции. Погрешность интерполирования	344
4.3.1. Введение	—
4.3.2. Интерполяционные формулы Лагранжа и Ньютона	347
4.3.3. Остаток интерполирования и его представления для некоторых классов функций	349

§ 4.4. Некоторые правила интерполирования при равноотстоящих значениях аргумента	358
4.4.1. Правила для интерполирования в начале и конце таблицы	—
4.4.2. Правила интерполирования внутри таблицы	360
§ 4.5. Приложение интерполирования к численному нахождению производных	364
4.5.1. Об интерполяционном правиле вычисления производной от функции, заданной таблично	—
4.5.2. Некоторые частные правила вычисления производных	369
§ 4.6. Интерполяционные методы решения численных уравнений	372
4.6.1. Введение. Связь с задачей обратного интерполирования	—
4.6.2. Метод приближений, основанный на интерполировании обратной функции	374
4.6.3. Замена точного уравнения $f(x)=0$ приближенным, полученным интерполированием f	376
§ 4.7. Интерполирование с кратными узлами	377
4.7.1. Существование и единственность интерполирующего многочлена. Остаток	—
4.7.2. Представление $R(x)$ в случае аналитической функции f . Формула Эрмита для многочлена $P(x)$	380
§ 4.8. Сходимость интерполяционных процессов	383
4.8.1. О предельной функции распределения узлов	384
4.8.2. Сходимость интерполирования аналитических функций	385
4.8.3. Некоторые вспомогательные теоремы	394
4.8.4. Сходимость интерполирования на множествах непрерывных и непрерывно дифференцируемых функций	399
Л и т е р а т у р а	413
ГЛАВА 5. ЧИСЛЕННОЕ ИНТЕГРИРОВАНИЕ	414
§ 5.1. Квадратурная сумма и условия ее построения. Остаток квадратуры	—
5.1.1. О квадратурной сумме	—
5.1.2. Остаток приближенной квадратуры	419
§ 5.2. Интерполяционные квадратурные правила и их погрешности	420
§ 5.3. Правила Ньютона — Котеса	424
§ 5.4. Некоторые простейшие правила Ньютона — Котеса	432
5.4.1. Правило трапеций	—
5.4.2. Правило парабол (Формула Симпсона)	434
5.4.3. Правило «трех восьмых»	436
§ 5.5. Квадратурные правила наивысшей алгебраической степени точности	437
5.5.1. Построение правила и его единственность	—
5.5.2. Два замечания о квадратурных коэффициентах	441
5.5.3. Остаток квадратурного правила	443
5.5.4. Сходимость квадратурного процесса наивысшей степени точности	444
5.5.5. Замечание об интегрировании периодических функций	446
§ 5.6. Некоторые частные случаи квадратурных правил наивысшей алгебраической степени точности	447
5.6.1. Постоянная весовая функция	—
5.6.2. Интегралы вида $\int_a^b (b-x)^{\alpha} (x-a)^{\beta} f(x) dx$	450
5.6.3. Интегралы вида $\int_0^{\infty} x^{\alpha} e^{-x} f(x) dx$	455
5.6.4. Интегралы вида $\int_{-\infty}^{\infty} e^{-x^2} f(x) dx$	457

§ 5.7. Квадратурные правила наивысшей степени точности, имеющие фиксированные заранее узлы	458
5.7.1. Некоторые общие теоремы	—
5.7.2. Некоторые частные квадратурные правила	461
§ 5.8. Квадратурные правила с равными коэффициентами	463
5.8.1. Построение формул Чебышева. Существование и единственность	—
5.8.2. Случай постоянного веса $p(x) \equiv 1$	466
§ 5.9. Увеличение точности квадратурных правил. Формулы эйлера вида	474
5.9.1. Введение	—
5.9.2. Правила эйлера вида	476
5.9.3. Формула Эйлера — Маклорена	480
5.9.4. Разностные видоизменения формулы Эйлера — Маклорена	485
§ 5.10. Увеличение точности квадратурных правил. Ослабление особенностей интегрируемой функции	487
§ 5.11. Сходимость квадратурного процесса	491
5.11.1. Условия сходимости общего квадратурного процесса	—
5.11.2. Сходимость интерполяционных квадратурных процессов	498
§ 5.12. Вычисление неопределенного интеграла	500
5.12.1. Введение	—
5.12.2. Погрешность вычислений и сходимость	504
§ 5.13. Понятие о некоторых частных методах вычисления неопределенного интеграла	514
5.13.1. Интегрирование функции, заданной таблицей значений	—
5.13.2. Вычисление при помощи периодически расположенных узлов	517
5.13.3. О правилах, использующих в вычислениях несколько предшествующих значений интеграла	521
Л и т е р а т у р а	531
ДОБАВЛЕНИЕ I. НЕКОТОРЫЕ СВЕДЕНИЯ ИЗ ФУНКЦИОНАЛЬНОГО АНАЛИЗА	532
§ 1. Метрические пространства. Сходимость и полнота	—
§ 2. Линейные нормированные пространства. Линейные операторы	535
§ 3. Дифференцирование нелинейных операторов и некоторые теоремы, с этим связанные	546
ДОБАВЛЕНИЕ II. ЧИСЛА И МНОГОЧЛЕНЫ БЕРНУЛЛИ	554
§ 1. Числа Бернулли	—
§ 2. Многочлены Бернулли и их свойства	556
§ 3. Периодические функции, связанные с многочленами Бернулли	560
§ 4. Представление произвольной функции при помощи многочленов Бернулли	562
ДОБАВЛЕНИЕ III. АЛГЕБРАИЧЕСКИЕ МНОГОЧЛЕНЫ НАИЛУЧШЕГО ПРИБЛИЖЕНИЯ	565
ДОБАВЛЕНИЕ IV. НЕКОТОРЫЕ СВЕДЕНИЯ ОБ УРАВНЕНИЯХ В КОНЕЧНЫХ РАЗНОСТЯХ	572
§ 1. Уравнения в конечных разностях произвольного вида	—
§ 2. Линейные уравнения	573
§ 3. Линейные уравнения с постоянными коэффициентами	578

