

ТБИЛИССКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ
ИНСТИТУТ ПРИКЛАДНОЙ МАТЕМАТИКИ
ИМЕНИ АКАДЕМИКА И.Н.ВЕКУА

С.Г.МИХЛИН

ПОГРЕШНОСТИ ВЫЧИСЛИТЕЛЬНЫХ ПРОЦЕССОВ

ИЗДАТЕЛЬСТВО ТБИЛИССКОГО УНИВЕРСИТЕТА
ТБИЛИСИ - 1983

6/8

|

Печатается по постановлению Ученого совета Института прикладной математики имени И.Н.Векуа Тбилисского государственного университета.

В книге рассматривается довольно широкий класс задач вычислительной математики и вычислительные процессы, которые приводят к построению решения данной задачи, анализируются источники погрешностей и оцениваются самые погрешности в предположении, что данная задача поставлена точно. Как кажется автору, в весьма многих случаях существует точно четыре источника погрешностей: погрешность аппроксимации, происходящая от замены данной задачи другой, упрощенной задачей; погрешность искажения, связанная с тем, что данные упрощенной задачи вычисляются неточно; погрешность алгоритма, возникающая в тех случаях, когда алгоритм, примененный для решения упрощенной задачи, дает после конечного числа шагов лишь приближенное её решение; погрешность округления, возникающая из-за того, что предписания упомянутого алгоритма выполняются неточно. Для широкого круга вычислительных процессов даны оценки всех перечисленных погрешностей.

Книга содержит 9 глав. В главе I дана постановка проблемы и исследованы для более простых примера. В главах II-VI исследуются погрешности линейных вычислительных процессов, в главах VII-IX тот же вопрос исследуется для нелинейных процессов.

Книга содержит расширенное изложение докладов, сделанных автором в Институте прикладной математики имени академика И.Н. Векуа в 1982 году.

Редактор: заслуженный деятель науки СССР

Л. Г. Магнарадзе

Рецензент: доктор физ.-матем. наук, проф.

В. Г. Мазья

ПРЕДИСЛОВИЕ

В настоящее время ни у кого не вызывает сомнения необходимость исследования и оценки различных погрешностей, возникающих при производстве вычислений, хотя ещё сравнительно недавно - лет 35-40 тому назад - многие видные представители прикладных наук считали, что в таком исследовании нет нужды и что любой вычислительный процесс надёжен, если он кажется таким. Следует признать, что в прошлом эта наименьшая точка зрения была не совсем безосновательной: до того, как в научный обиход вошли современные вычислительные машины, доступными для реализации были, как правило, только методы вычислений, дававшие грубые приближения; такие методы требовали небольшого объёма вычислений, при котором не могло произойти катастрофических вычислительных погрешностей. Положение изменилось, когда электронные вычислительные машины получили широкое распространение. Стало возможным выполнение вычислений огромного объёма, но при этом довольно скоро обнаружилось, что в ряде случаев результат вычислений становится совершенно недоверенным из-за накопления погрешностей счета. Специалисты по вычислительной математике в своё время обнаружили это явление и исследовали частный, но весьма важный вопрос о накоплении погрешностей округления при решении линейных алгебраических систем /Уилкинсон и др./. В своих работах, относящихся к началу 60-х годов, автор этих строк обратил внимание на погрешности /"погрешности искажения"/ возникающие от неточной замены данной задачи на более простую приближённую задачу. С погрешностью искажения связано введённое тогда же автором понятие устойчивости вычислительного процесса. В последние годы автор занимался выявлением всех возможных источников погрешностей и оценкой этих погрешностей; предлагаемая вниманию читателей книга в основном содержит результаты этих исследований. Наряду с процессами, которые автор называет свободными /в них результат любого шага не зависит от результатов остальных шагов процесса/, исследованы также рекуррентные процессы. Много внимания уделено нелинейным вычислительным процессам, в частности, тем, которые возникают при решении односторонних вариационных задач.

Ниже следует более подробное изложение содержания книги. Она содержит 9 глав. В главе I описан класс задач вычислительной математики, которые решаются так, что каждая задача заменяется последовательностью более простых задач /точнее, любой задачей из этой последовательности/. Такая замена приводит к тому, что в данной книге названо погрешностью аппроксимации - это некоторый функционал, ограничивающий погрешность от замены данной задачи на более простую. Далее, данные приближенной задачи на самом деле нам не заданы - они нами вычисляются и притом, за редкими исключениями, с некоторой погрешностью. Таким образом, мы получаем приближенную задачу в искаженном виде; если искаженную задачу решить точно, то мы придем не к точному, а к искаженному решению приближенной задачи. Всё это приводит нас ко второму виду погрешности, которую автор называет погрешностью искажения. Как теоретические рассуждения, так и численные примеры показывают, что погрешность искажения для решения может быть очень велика, даже если погрешности данных приближенной задачи малы.

Допустим, что мы умеем решать приближенную задачу, независимо от того, дана она нам в точном или в искаженном виде. Это значит, что нам известен некоторый алгоритм, на каждом шаге которого требуется выполнить конечное число операций, выполнение которых нам доступно, и после некоторого конечного числа шагов этот алгоритм приводит /если упомянутые операции выполнены точно, без погрешности/ либо к точному решению искаженной приближенной задачи, либо к элементу, который в том или ином смысле близок к названному решению, причем эту близость можно оценить. Таким путем мы приходим к третьему виду погрешности - к погрешности алгоритма.

Последняя - четвертая - погрешность есть погрешность округления; она связана с тем, что предписания упомянутого выше алгоритма выполняются не точно, а с ошибками округления.

В 10^й же главе I рассмотрены два примера: квадратурные формулы и системы линейных алгебраических уравнений с квадратной матрицей и с неравным нулю определителем. Первый пример интересен тем, что для него существенна только погрешность аппроксимации; во втором примере эта погрешность отсутствует, хотя, например, по методу Гаусса исходная задача заменяется другой, бо-

лее простой. Вообще, в каждом конкретном случае та или другая из отмеченных выше четырех погрешностей может отсутствовать.

Главы II-VI посвящены погрешностям линейных вычислительных процессов. В главе II для ряда приближенных методов изучается погрешность аппроксимации, в главе III - погрешность искажения, в главе IV - погрешности алгоритма и округления. Глава V содержит анализ погрешностей для метода конечных элементов, глава VI - для различных методов решения интегральных уравнений, Fredholmовских и сингулярных.

В главах VII-IX изучаются погрешности нелинейных вычислительных процессов. В главе VII рассмотрены общие вопросы; полученные результаты применяются к методу Рунге /в частности, к методу конечных элементов/ для нелинейных задач. Подробно исследованы нелинейные рекуррентные процессы, включая метод Ньютона - Канторовича.

В главе VIII рассмотрены приближенные решения одного класса односторонних вариационных задач. Исследованы погрешности аппроксимации и искажения; доказана устойчивость точного решения относительно малых возмущений данных; даны оценки возмущения решения в зависимости от возмущений данных задачи.

Глава IX содержит приложения результатов главы VIII к задачам упруго-пластического состояния, основанные на теории Сен-Венана - Мизеса и вариационном принципе Хаара - Кармана. В § 1 приведены, для удобства ссылок, основные факты теории Сен-Венана - Мизеса и принцип Хаара - Кармана. В § 2 рассмотрена задача упруго-пластического кручения. Отмечены существование решения и эквивалентность принципа Хаара - Кармана системе уравнений упруго-пластического кручения /все эти результаты получены французскими математиками школы Лиюкса/. Далее, доказана устойчивость точного решения задачи относительно малых возмущений данных задачи - крутящего момента и поперечного сечения стержня; получены оценки погрешностей аппроксимации и искажения. В § 3 аналогичные результаты получены для плоской задачи упруго-пластического состояния; в частности, при некоторых предположениях доказана эквивалентность принципа Хаара - Кармана и уравнений плоской задачи. Аналогичные результаты получены в § 4 для трехмерной задачи.

И.Н.Гончарова провела все вычисления по одному из примеров § 4 главы III. В.Б.Тихтин написал добавление к главе УШ, в котором дан обзор ряда современных работ, содержащих оценки погрешности аппроксимации для решений довольно широкого класса вариационных неравенств. Автор рад выразить И.Н.Гончаровой и В.Б.Тихтину свои сердечную благодарность.

Июнь 1982 г.
Ленинград

С.Г.Михлин

ГЛАВА I

Постановка проблемы

§ I. Задачи вычислительной математики. Вычислительные процессы

1°. Весьма многие задачи вычислительной математики можно подвести под следующую схему. Дан некоторый объект f ; требуется определить другой объект x , удовлетворяющий двум требованиям: 1) x принадлежит некоторому данному классу объектов X ; 2) x находится в некотором данном отношении \mathcal{U} с данным объектом f . Нашу задачу можно записать так:

$$\{x \in X, x \mathcal{U} f\}. \quad (I.I.1)$$

Объект x , удовлетворяющий условиям (I.I.1), есть решение нашей задачи.

Примеры. I. Пусть требуется вычислить интеграл

$$\int_0^1 f(t) dt,$$

где f - заданная функция. В этом случае $f = f(t)$, $X = R_1$, а соотношение \mathcal{U} означает, что x и f связаны уравнением

$$x = \int_0^1 f(t) dt. \quad (I.I.2)$$

В несколько усложненном примере, когда требуется вычислить интеграл вида

$$x(x) = \int_{\mathcal{D}} f(x, t) dt, \quad (I.I.3)$$

где \mathcal{D} - множество в R_m , а x - параметр, который может принять любое значение из некоторого множества Ω , f есть функция двух точек x и t , X есть множество всех функций, определенных на Ω , а соотношение \mathcal{U} определяется равенством (I.I.3).

2°. Пусть требуется решить систему линейных алгебраических уравнений

$$Ax = b, \quad (I.I.4)$$

где x и b - k -мерные векторы, A - матрица порядка $k \times k$. Пусть для простоты все числа в нашей задаче - действительные. Тогда $f = b$, $X = R_k$, соотношение \mathcal{U} определяется уравнением (I.I.4).

3. Рассмотрим смешанную задачу для волнового уравнения:

$$\frac{\partial^2 x}{\partial t^2} - \Delta x = f(x, t), \quad t \geq 0, \quad x \in \Omega, \quad (I.I.5)$$

$$x(x, 0) = \varphi(x), \quad x_t(x, 0) = \psi(x),$$

$$x(x, t)|_{x \in \partial \Omega} = \omega(x, t);$$

Δ - оператор Лапласа. Будем, например, искать решение этой задачи, принадлежащее соболевскому классу $W_p^{(2)}([0, \infty) \times \Omega)$. Тогда $X = W_p^{(2)}([0, \infty) \times \Omega)$, f есть совокупность функций $f(x, t)$, $\varphi(x)$, $\psi(x)$, $\omega(x, t)$, а соотношение τ определяется совокупностью уравнений (I.I.5).

4. Пусть $J_f(X)$ - функционал, определенный на некотором множестве X , и требуется найти элемент $x \in X$, на котором $f(x)$ достигает минимума. В данном случае $f = J_f(x)$, а соотношение τ можно записать так: если x_* - решение нашей вариационной задачи, то $x_* \tau_f$ означает, что $f(x_*) \leq f(x), \forall x \in X$, или, если угодно, что $f(x_*) = \min_{x \in X} f(x)$.

3°. Как правило, задача (I.I.I) решается приближенно. В довольно широком классе случаев это означает следующее. Каждому натуральному n приводится в соответствие: объект $X_n^{(n)}$; класс объектов X_n ; оператор P_n , отображающий X_n в X ; соотношение τ_n . Задача (I.I.I) заменяется любой задачей из последовательности

$$\{U^{(n)} \in X_n, U^{(n)} \tau_n f^{(n)}\}; \quad (I.I.6)$$

мы предполагаем, что решение каждой задачи этой последовательности не зависит от решений последующих задач, но может зависеть от предыдущих. Если $U_*^{(n)}$ есть решение n -й задачи (I.I.6), то элемент $x_* = P_n U_*^{(n)}$ рассматривается как приближенное решение задачи (I.I.I).

Будем говорить, что задачи (I.I.6) аппроксимируют задачу (I.I.I) и допустим, что аппроксимирующие задачи (I.I.6) в том или ином смысле проще, чем аппроксимируемая задача (I.I.I). Процесс решения последовательности задач (I.I.6) назовем вычислительным процессом. Этот процесс будем называть свободным, если при любом n решение $U_*^{(n)}$ определяется независимо от значений элементов $U_*^{(k)}, k \neq n$, и рекуррентным, если $U_*^{(k)}$ определяется

через $U_*^{(k)}$, $k < n$.

Пример 5. Будем вычислять интеграл (I.1.2) (пример I), например, по формуле средних прямоугольников с n узлами. В таком случае $f^{(n)}$ есть совокупность n ординат $f(t_k)$ в точках $t_k = \frac{2k-1}{2n}$, $k = 1, 2, \dots, n$; $X_n = R_1$, R_n есть тождественный оператор в R_1 , а соотношение U_n определяется выбранной квадратурной формулой:

$$U_n^{(n)} U_n f^{(n)} := U^{(n)} = \frac{1}{n} \sum_{k=1}^n f\left(\frac{2k-1}{2n}\right). \quad (I.1.7)$$

Вычислительный процесс (I.1.7) - свободный.

Пример 6. В примере 2 положим $A = I - B$, где I - единичная матрица в R_k , и запишем уравнение (I.1.4) в виде $x = Bx + b$. Воспользуемся методом простой итерации: зададим какой-нибудь вектор $U^{(0)}$ и положим для любого натурального n

$$U^{(n)} = B U^{(n-1)} + b. \quad (I.1.8)$$

В данном случае $f^{(n)}$ есть совокупность векторов b и $U^{(n-1)}$, $X_n = X = R_k$, $R_n = I$, соотношение U_n определяется формулой простых итераций

$$U_n^{(n)} U_n f^{(n)} := U^{(n)} = B U^{(n-1)} + b. \quad (I.1.9)$$

Вычислительный процесс (I.1.9) - рекуррентный.

Приведем еще некоторые примеры. К свободным вычислительным процессам приводят методы Рунге, Лунгова - Галеркина, конечных элементов, разностные методы, метод прямых. К рекуррентным вычислительным процессам приводят различные итерационные процессы: метод простой итерации, релаксационные методы, градиентные методы, метод наискорейшего спуска, метод Ньютона - Канторовича и др.

4^o. Предположим, что аппроксимирующую задачу (I.1.6) мы "умеем решать" при любом n . Под этим мы понимаем следующее. Существует набор "элементарных" операций, которые мы умеем выполнять, и мы располагаем алгоритмом, который при любом n после выполнения конечного числа "элементарных" операций вычислит (если эти операции выполнены точно, без округности) либо к точному значению элемента $U_*^{(n)}$, решающего задачу (I.1.6), либо к

некоторому элементу $\tilde{U}^{(n)}$, который в том или ином смысле близок к элементу $U_*^{(n)}$, причем степень этой близости можно оценить. Разумеется, мы предполагаем при этом, что для рассматриваемых n аппроксимирующая задача разрешима.

Поясним сказанное на примере. Пусть задаче (I.I.I) имеет вид

$$Ax = f, \quad (\text{I.I.IO})$$

где A - положительно определенный оператор в некотором гильбертовом пространстве H , x и f - элементы этого пространства, искомый и данный. Если к этой задаче применить какой-нибудь вариант метода Рунца (например, метод конечных элементов), то мы придем к аппроксимирующей линейной алгебраической задаче вида

$$A_n u^{(n)} = f^{(n)}, \quad (\text{I.I.II})$$

где A_n - числовая матрица порядка $N \times N$, $u^{(n)}$ и $f^{(n)}$ - N -мерные числовые векторы; N - некоторая целочисленная функция от n . Эта задача разрешима при любом n . Если для ее решения воспользоваться алгоритмом Гаусса, то через конечное число шагов мы придем к точному решению $U_*^{(n)}$; если же воспользоваться, например, алгоритмом простой итерации (для этого задачу (I.I.II) надо предварительно подвергнуть некоторому элементарному преобразованию), то в результате конечного числа шагов мы получим некоторый вектор $\tilde{U}^{(n)}$, близкий к вектору $U_*^{(n)}$. Набор элементарных операций, о котором говорилось выше, в обоих случаях содержит арифметические действия; в случае метода Гаусса в этот набор входит еще определение наибольшего числа из заданного конечного множества чисел.

5°. Для приложений наиболее интересен случай, когда X и X_n суть банаховы пространства. Предположим, еще, что f и $f^{(n)}$ также суть элементы некоторых банаховых пространств Y и Y_n соответственно. Эти предположения мы сохраним на протяжении всей книги.

§ 2. Источники погрешностей.

1°. Если задача (I.I.I) решается каким-нибудь методом, подходящим под схему § I, то, как кажется автору, речь может идти о четырех возможных источниках погрешностей. При этом мы считаем, что задача (I.I.I) поставлена точно. Тем самым мы пренебрегаем возможными погрешностями, связанными с постановкой упомянутой задачи как задачи естествознания или социально-научных

дисциплины (погрешности измерений, недостаточно точности основных гипотез и т.п.). По мнению автора, устранение или хотя бы уменьшение подобных погрешностей лежат вне пределов действия математических методов, и самое большее, на что здесь можно рассчитывать, это, при наличии достаточной информации о степени точности упомянутых выше факторов, воспользоваться методами теории возмущений и дать оценку погрешности, обусловленной неточностью этих факторов.

2°. Переходим к описанию источников погрешностей.

1. Погрешность аппроксимации. Переход от точно поставленной задачи (I.I.I) к аппроксимирующей задаче (I.I.6) с выбранным номером n приводит к тому, что элемент X_* пространства X заменяется элементом $X_*^{(n)} = \rho_n U_*^{(n)}$ того же пространства. Возникающую при этом погрешность $X_* - \rho_n U_*^{(n)}$ естественно оценить ее нормой

$$\rho_n = \|X_* - \rho_n U_*^{(n)}\| = \|X_* - X_*^{(n)}\|. \quad (I.2.I)$$

Величину ρ_n будем называть погрешностью аппроксимации задачи (I.I.I).

Вопрос об оценке погрешности аппроксимации - один из важнейших в вычислительной математике, а ему посвящена обширная литература.

2. Погрешность искажения. Как правило, переход от точной задачи (I.I.I) к аппроксимирующей задаче (I.I.6) совершается с погрешностью. Элементы $\rho^{(n)}$ не заданы заранее, а вычисляются по данным точной задачи. Вычислительные операции почти всегда выполняются с той или иной погрешностью, и вместо элементов $\rho^{(n)}$ мы на самом деле получаем некоторые "искаженные" элементы $\tilde{\rho}^{(n)}$. Может случиться, что и соотношения U_n также заменяются "искаженными" соотношениями \tilde{U}_n . Элементы $\rho^{(n)}$ и соотношения U_n на самом деле остаются неизвестными; известными являются $\tilde{\rho}^{(n)}$ и \tilde{U}_n . В результате вместо последовательности аппроксимирующих задач (I.I.6) получится последовательность "искаженных аппроксимирующих задач"

$$n \geq 1. \{z \in X_n, z \sim \tilde{U}_n \tilde{\rho}^{(n)}\}. \quad (I.2.2)$$

Допустим, что эти задачи разрешимы, а пусть z_* их решение. Величину

$$\hat{\xi}_n = \| \hat{z}_x^{(n)} - U_x^{(n)} \|_{X_n} \quad (1.2.3)$$

назовем " X_n -погрешностью искажения". Большой интерес представляет величина

$$\xi_n = \| \rho_n z_x^{(n)} - \rho_n U_x^{(n)} \|_X, \quad (1.2.4)$$

которую мы назовем X -погрешностью, или просто погрешностью искажения.

3. Искаженную задачу (1.2.2) предстоит решить. Будем считать (см. § I), что мы располагаем некоторым алгоритмом, который после выполнения конечного числа действий приводит, если эти действия выполнены без погрешности, либо к точному значению элемента $z_x^{(n)}$, либо к такому элементу $w^{(n)} \in X_n$, что погрешность

$$\eta_n = \| \rho_n w^{(n)} - \rho_n z_x^{(n)} \|_X \quad (1.2.5)$$

можно оценить и можно добиться того, чтобы величина, оценивающая эту погрешность, была сколь угодно малой. Первый случай - частный по отношению ко второму, и мы будем рассматривать именно этот второй случай, допуская, что величина (1.2.5) может иногда равняться нулю. Эту величину мы назовем X -погрешностью, или просто погрешностью алгоритма; можно ввести и " X_n -погрешность алгоритма"

$$\hat{\eta}_n = \| w^{(n)} - z_x^{(n)} \|_{X_n}. \quad (1.2.6)$$

Отметим, что погрешность алгоритма можно рассматривать как погрешность аппроксимации для задачи (1.2.2).

4. Практически любой алгоритм содержит ряд операций, выполнение которых сопряжено с погрешностями. Таковы арифметические действия, операции вычисления значений функций по таблицам или по приближенным формулам и т.п. Названные погрешности приводят к тому, что вместо элемента $w^{(n)}$ мы получим другой элемент $\tilde{w}^{(n)} \in X_n$. Величину

$$\hat{\zeta}_n = \| \tilde{w}^{(n)} - w^{(n)} \|_{X_n} \quad (1.2.7)$$

назовем " X_n -погрешностью округления", а величину

$$\zeta_n = \| P_n \tilde{w}^{(n)} - P_n w^{(n)} \|_X \quad (1.2.8)$$

- " X -погрешностью округления", или просто "погрешностью округления"

3°. В конечном счете вместо искомого элемента $\alpha_x \in X$ мы вычислим элемент $\tilde{w}^{(n)} \in X_n$. За приближенное решение задачи (1.1.1) естественно принять элемент $P_n \tilde{w}^{(n)}$; его погрешность

$$\delta_n = \| \alpha_x - P_n \tilde{w}^{(n)} \| \leq \rho_n + \xi_n + \eta_n + \zeta_n. \quad (1.2.9)$$

4°. Приведем некоторые литературные указания. Вопросы оценки погрешности аппроксимации возникли почти одновременно с возникновением анализа. Достаточно указать, например, на оценки остаточного члена степенного ряда. Из многочисленных работ более близких лет, в которых изучается погрешность аппроксимации, мы отметим здесь монографии Вязова и Форсайта [1], Вайлиско [10], Варга [1], Габдулхаева [3], Гавурина [1], Канторогича и Акилова [1], Канторовича и Крылова [1], Красносельского и др. [1], Марчука [1], автора этой книги [3, 22], Обана [1], Станесана и Руховца [1], Пресдорфа и Зильбермана [1], Соболева [2], Странга и Фикса [1], Сьярле [1]. Понятие погрешности искажения было введено первоначально в работах автора [4, 5, 7]; в первых двух исследована устойчивость классического метода Рунге по отношению к погрешностям искажения, в третьей получены некоторые общие теоремы об устойчивости вычислительных процессов, которые выше были названы свободными. Самый термин "искажение" ввела Л.Н. Довбыт в работе [1], посвященной погрешностям искажения собственных чисел и собственных подпространств самосопряженных операторов. В работе Ясковой и Яковлева [1] исследована устойчивость относительно погрешностей искажения для метода Бубнова - Галеркина. Велиев [1 - 9] исследовал устойчивость метода Бубнова - Галеркина для ряда нестационарных задач, линейных и нелинейных. Перечисленные здесь работы, опубликованные не позднее 1965 г., суммированы в книге автора [8], в которой кроме того содержатся некоторые результаты того же рода, но относящиеся к нелинейным задачам. Дальнейшие результаты по погрешностям искажения свободных вычислительных процессов, линейных и нелинейных, содержатся в работах автора [12, 22].

Яношича [1], Таккера [1], Омодеи [1], Суворова [1]. Тот же вопрос для рекуррентных вычислительных процессов рассмотрен в работах автора [33, 34].

Вопрос об оценке погрешности алгоритма в большинстве случаев достаточно прост и, насколько известно автору, специально не освещался.

Много работ посвящено погрешностям округления; эти работы относятся главным образом к различным методам решения линейных алгебраических систем. Отметим здесь работы Голуба [1], Уилкинсона [1, 2] и Воеводина [1]. Некоторые другие задачи рассмотрены автором [33, 34].

В заключение отметим, что некоторая классификация погрешностей предложена в учебнике Березина и Жидкова [1]. Эти авторы различают три типа погрешностей: 1) неустранимая погрешность, вытекающая из неточностей постановки задачи (1.1.1); 2) погрешность метода - она совпадает с определенной выше погрешностью аппроксимации; 3) вычислительная погрешность - она включает в себя введенные выше погрешности искажения, алгоритма и округления.

Сображения автора относительно неустранимой погрешности изложены (без применения самого термина) выше, в п.1⁰ настоящего параграфа. Что касается вычислительной погрешности, то нам хотелось бы привести некоторые доводы в пользу нашей более подробной классификации. Погрешность искажений зависит лишь от того, с какой точностью вычислены объекты, которые являются данными в приближенной задаче (1.2.1); она не зависит от того, какой алгоритм применен для решения этой последней задачи и с какой точностью выполнены предписания этого алгоритма. От его выбора зависят погрешность алгоритма. Наконец, погрешность округления полностью определяется при выбранном алгоритме той точностью, с которой выполняются предписания алгоритма. Таким образом, вычислительная погрешность Березина и Жидкова распадается на три слагаемых, каждое из которых имеет свой источник и оценивается независимо.

Наша классификация дана в работе автора [33].

§. Погрешности квадратурных формул.

Чтобы проиллюстрировать общие соображения § 2, мы проанализируем в следующем параграфе погрешность обычных квадратурных

формул.

Погрешность аппроксимации для этих формул изучалась с давних времен. Мы приведем здесь оценки этой погрешности для наиболее распространенных квадратурных формул. Через (a, b) обозначен интервал интегрирования, через n - число промежутков разбиения, через f - подынтегральная функция.

1) формула прямоугольников

$$\rho_n \leq \frac{(b-a)^2}{2n} \|f'\|_{C[a,b]};$$

2) формула средних прямоугольников

$$\rho_n \leq \frac{(b-a)^3}{24n^2} \|f''\|_{C[a,b]};$$

3) формула трапеций

$$\rho_n \leq \frac{(b-a)^3}{12n^2} \|f''\|_{C[a,b]};$$

4) формула Симпсона

$$\rho_n \leq \frac{(b-a)^5}{180n^4} \|f^{(4)}\|_{C[a,b]};$$

5) формула Гаусса, $(a, b) = (-1, 1)$,

$$\rho_n \leq \frac{2^{2n+1} (n!)^4}{(2n+1) [(2n)!]^3} \|f^{(2n)}\|_{C[-1, 1]};$$

6) формула Эйлера - Маклорена

$$\rho_n \leq \frac{nh^{2m+1} B_{2m}}{(2m)!} \|f^{(2m)}\|_{C[a,b]};$$

в последней формуле m - произвольно выбранное натуральное число, от которого зависит точность формулы Эйлера - Маклорена, $h = (b-a)/n$, B_{2m} - числа Бернулли.

Отметим еще, что в книге Соболева [2] построены кубатурные формулы, которые для подынтегральных функций класса $W_2^{(m)}(\Omega)$ (Ω - область интегрирования) приводят к наименьшей или близкой к наименьшей погрешности аппроксимации.

Исследование остальных погрешностей мы проведем для класса квадратурных формул вида

$$\int_a^b f(t) dt \approx \sum_{k=1}^n c_k f(t_k), \quad (1.3.1)$$

обладающих следующими свойствами: $C_k \geq 0$; формула (1.3.1) точна, если $f(t) = \text{const}$, так что

$$\sum_{k=1}^n C_k = b-a. \quad (1.3.2)$$

На узлы t_k мы не накладываем никаких ограничений. Подынтегральную функцию предполагаем ограниченной: $|f(t)| \leq M = \text{const}$. Пусть коэффициенты C_k и ординаты $y_k = f(t_k)$ вычислены точно, так что на самом деле нам известны их искаженные значения $\tilde{C}_k = C_k + \gamma_k$ и $\tilde{y}_k = y_k + d_k$. Мы предполагаем при этом, что

$$|d_k| \leq \delta, \quad \sum_{k=1}^n |\gamma_k| \leq \gamma, \quad (1.3.3)$$

где δ и γ - малые числа, и что искаженные коэффициенты \tilde{C}_k удовлетворяют условиям, сформулированным выше:

$$\tilde{C}_k \geq 0, \quad \sum_{k=1}^n \tilde{C}_k = b-a.$$

Погрешность искажения формулы (1.3.1) равна

$$\begin{aligned} \epsilon_n &= \left| \sum_{k=1}^n [(C_k + \gamma_k)(y_k + d_k) - C_k y_k] \right| = \\ &= \left| \sum_{k=1}^n (\gamma_k y_k + \tilde{C}_k d_k) \right| \leq M\gamma + \delta. \end{aligned}$$

Таким образом, малые искажения коэффициентов и ординат приводят к малому искажению интеграла, вычисленного по квадратурной формуле.

После того, как ординаты вычислены, алгоритм вычисления правой части формулы (1.3.1) содержит только конечное число действий сложения и умножения, поэтому для квадратурной формулы погрешность алгоритма равна нулю.

Оценим, наконец, погрешность округления. Интеграл $\int_a^b f(t) dt$ мы вычисляем по искаженной квадратурной формуле

$$\int_a^b f(t) dt \approx \sum_{k=1}^n \tilde{C}_k \tilde{y}_k. \quad (1.3.4)$$

Предположим, что каждое из произведений $\tilde{C}_k \tilde{y}_k$ вычисляется с погрешностью, не превосходящей числа δ^p , где δ - основание системы счисления, p - некоторое положительное число. Тогда погрешность округления формулы (1.3.4) не превосходит числа $n\delta^p$.

Если желательнее, чтобы эта погрешность была не больше заданного числа ϵ , то достаточно подчинить ρ неравенству

$$\rho \geq \frac{\ln(n/\epsilon)}{\ln 5}. \quad (I.3.5)$$

Нетрудно проанализировать также погрешности квадратурных формул типа формулы Эйлера - Маклорена, отличающихся наличием некоторого фиксированного числа отрицательных коэффициентов. Не вносит новых трудностей и исследование погрешностей кубатурных формул.

§ 4. Погрешности решения системы линейных алгебраических уравнений.

1°. Пусть дана система n линейных алгебраических уравнений о n неизвестными, которая в матричной форме записывается в виде

$$Ax = f. \quad (I.4.1)$$

Искомый вектор x и данный вектор f будем рассматривать как элементы пространства R_n , а матрицу A -- как оператор, действующий в том же пространстве. Допустим, что этот оператор обратим, так что $\|A^{-1}\| < \infty$.

Будем решать систему (I.4.1) по методу Гаусса с применением схемы единственного деления по главным элементам (см. Фаддеев и Фаддеева [1], § 16). Метод Гаусса состоит из прямого и обратного хода. Прямой ход данная система преобразуется в систему с треугольной матрицей, главная диагональ которой состоит из единиц; обратный ход заключается в решении новой, рекуррентной системы. С точки зрения понятий, изложенных в § 1, прямой ход можно рассматривать как преобразование задачи (I.1.1) в задачу (I.1.6); при этом $X_n = X = R_n$, $x^{(n)} = x$, $axf := Ax = f$, $x^{(n)} x_n f^{(n)} := Lx = f$, где L - упомянутая выше треугольная матрица, а f - столбец новых свободных членов, полученных в результате прямого хода. В данном случае погрешность аппроксимации равна $1/n!$, но возникает погрешность искажения. Исследуем ее.

Действия, описанные выше, требуют, чтобы делителями были отличны от нуля, а это, в свою очередь, требует перенумерации уравнений и неизвестных. Будем считать, что перед каждым уравнением такая перенумерация выполнена, и не будем этого особо отмечать.

вать.

Впишем первое уравнение системы

$$a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n = f_1.$$

Разделим его на a_{11} . Если бы деление было выполнено точно, мы получили бы уравнение

$$x_1 + b_{12}^{(1)} x_2 + \dots + b_{1n}^{(1)} x_n = \tilde{f}_1^{(1)}, \quad (1.4.2)$$

где $b_{1k}^{(1)} = a_{1k}/a_{11}$, $\tilde{f}_1^{(1)} = f_1/a_{11}$. Однако, в общем случае деление совершается с погрешностью, и мы приходим к уравнению вида

$$x_1 + (b_{12}^{(1)} + r_{12}) x_2 + \dots + (b_{1n}^{(1)} + r_{1n}) x_n = \tilde{f}_1 + \delta_1.$$

Следующий шаг заключается в исключении неизвестной x_1 из остальных $(n-1)$ уравнений. Это приведет нас к системе уравнений

$$\sum_{k=2}^n (b_{jk}^{(2)} + r_{jk}^{(2)}) x_k = \tilde{f}_j^{(2)} + \delta_j; \quad 2 \leq j \leq n. \quad (1.4.3)$$

Здесь через $b_{jk}^{(2)}$ и $\tilde{f}_j^{(2)}$ обозначены числа, которые получились бы, если бы все описанные выше действия были выполнены точно. Продолжая тот же процесс исключения неизвестных, получим систему вида

$$\sum_{k=j}^n (b_{jk} + r_{jk}) x_k = \tilde{f}_j + \delta_j; \quad 1 \leq j \leq n$$

или, короче

$$(L + \Gamma) x = \tilde{f} + \delta. \quad (1.4.4)$$

В системе (1.4.4) неизвестный вектор обозначен через x вместо X , так как из-за погрешностей Γ и δ системы (1.4.1) и (1.4.4) не равносильны, вообще говоря, и их решения, как правило, не совпадают.

Оценим погрешность $\|x - \alpha\|_{R_n} := \|x - \alpha\|$

$$\begin{aligned} \|x - \alpha\| &= \|(L + \Gamma)^{-1} (\tilde{f} + \delta) - L^{-1} \tilde{f}\| \leq \\ &\leq \|(L + \Gamma)^{-1} \tilde{f} - L^{-1} \tilde{f}\| + \|(L + \Gamma)^{-1} \delta\|. \end{aligned}$$

Первый член справа оценивается так: если $\|L^{-1}\| \|\Gamma\| \leq \beta < 1$, то

$$\|(L + \Gamma)^{-1} \tilde{f} - L^{-1} \tilde{f}\| \leq \|(L + \Gamma)^{-1} - L^{-1}\| \|\tilde{f}\| \leq \frac{\|L^{-1}\| \|\Gamma\| \|\alpha\|}{1 - \beta}.$$

оценку

$$|\gamma_k| \leq C_0 \max_{j \geq k+1} |d_j| \leq C_0 \varepsilon ,$$

в которой C_0 есть норма обратной матрицы системы (1.4.6), рассматриваемой как оператор в пространстве n -мерных векторов, в котором норма определяется по формуле $\|A\| = \max |a_{ij}|$. Таким образом, если $\gamma = (\gamma_1, \gamma_2, \dots, \gamma_n)$ есть вектор погрешности округления, то

$$\|\gamma_n\|_{R_n} \leq C_0 \sqrt{n} \varepsilon . \quad (1.4.8)$$

З^о. Пусть $A = I - B$ и $\|B\|_{R_n \rightarrow R_n} = q < 1$. Тогда уравнение (1.4.1) можно решать итерациями по формуле

$$x^{(k)} = Bx^{(k-1)} + f, \quad k = 1, 2, \dots ; \quad (1.4.9)$$

начальное приближение $x^{(0)}$ выбирается по произволу. Примем, что $x^{(0)} = 0$. В данном случае погрешности аппроксимации и искажения отсутствуют; погрешность алгоритма, если выполнены N итераций, равна

$$\|x_* - x^{(N)}\|_{R_n} = \left\| \sum_{k=N+1}^{\infty} B^k f \right\|_{R_n} \leq \frac{q^{N+1}}{1-q} \|f\|_{R_n} .$$

Оценим погрешность округления. Как и в п.2^о, ошибки округления приведут к тому, что вместо векторов $x^{(k)}$ мы получим некоторые другие векторы $w^{(k)}$. Обозначим через d_k результат ошибок округления на k -м шаге, так что

$$w^{(k)} = f + Bw^{(k-1)} + d_k . \quad (1.4.10)$$

Допустим на этот раз, что вычисления производятся с заданной относительной погрешностью, так что $\|d_k\|_{R_n} \leq \varepsilon \|w^{(k)}\|_{R_n}$. Точные значения векторов $x^{(k)}$ суть

$$x^{(k)} = f + Bx^{(k-1)} .$$

Фактически мы получим последовательность векторов

$$w^{(0)} = 0, \quad w^{(1)} = f, \quad w^{(2)} = f + Bx^{(1)} + d_2 = x^{(2)} + d_2$$

и т.д. Положим $\gamma_0 = \gamma_1 = 0$, $\gamma_2 = d_2$. Далее, $w^{(3)} = f + Bw^{(2)} = f + Bx^{(2)} + Bd_2 + d_3 = x^{(3)} + B\gamma_2 + d_3 := x^{(3)} + \gamma_3$.

Пусть для некоторого K оказалось $w^{(K-1)} = x^{(K-1)} + r_{K-1}$.

Тогда

$$w^{(K)} = f + Bw^{(K-1)} + d_K = f + Bx^{(K-1)} + Br_{K-1} + d_K,$$

что можно представить в виде $w^{(K)} = x^{(K)} + r_K$, где

$$r_K = Br_{K-1} + d_K. \quad (1.4.11)$$

Отсюда

$$\|r_K\| \leq q \|r_{K-1}\| + \|d_K\| \leq q \|r_{K-1}\| + \varepsilon \|r_K\| + \varepsilon \|x^{(K)}\|. \quad (1.4.12)$$

Но

$$\|x^{(K)}\| = \|f + Bf + \dots + B^{K-1}f\| \leq \frac{1}{1-q} \|f\|,$$

и из (1.4.12) находим

$$\|r_K\| \leq q_1 \|r_{K-1}\| + \varepsilon_1; \quad q_1 = \frac{q}{1-\varepsilon}, \quad \varepsilon_1 = \frac{\varepsilon \|f\|}{(1-\varepsilon)(1-q)}.$$

Изно, что $\|r_K\| \leq \tau_K$, где τ_K есть решение разностного уравнения

$\tau_K = q_1 \tau_{K-1} + \varepsilon_1$, $\tau_0 = 0$; это решение есть

$$\tau_K = \frac{\varepsilon_1}{1-q_1} (1-q_1^K)$$

и погрешность округления после N итераций имеет оценку

$$\|r_N\| \leq \frac{\varepsilon_1}{1-q_1} (1-q_1^N) < \frac{\varepsilon \|f\|}{(1-q_1)(1-q)(1-\varepsilon)}. \quad (1.4.13)$$

4°. Выше мы предполагали, что матрица A и столбец f заданы точно. Однако в практике применения приближенных методов мы чаще всего встречаемся с такой ситуацией, когда указанные величины известны с некоторыми погрешностями, так что решать приходится искаженную систему вида

$$(A + \Gamma) x = f + \delta.$$

Мы вправе считать, что нормы $\|\Gamma\|$ и $\|\delta\|$ достаточно малы; в частности, будем считать, что $\nu = \|A^{-1}\| \cdot \|\Gamma\| < 1$. Оценим погрешность искажения $\rho = \|x - x_{\text{идеал}}\|_{R_n}$. Имеем

$$z - \alpha = [(A + \Gamma)^{-1} - A^{-1}]f + (A + \Gamma)^{-1}\delta.$$

Но $(A + \Gamma)^{-1} = (I + A^{-1}\Gamma)^{-1}A^{-1}$ и $(A + \Gamma)^{-1} - A^{-1} = -(I + A^{-1}\Gamma)^{-1}A^{-1}\Gamma A^{-1}$, поэтому

$$z - \alpha = -(I + A^{-1}\Gamma)^{-1}A^{-1}\Gamma\alpha + (I + A^{-1}\Gamma)^{-1}A^{-1}\delta.$$

Далее, $\|(I + A^{-1}\Gamma)^{-1}\| \leq (1 - \nu)^{-1}$ и

$$\|z - \alpha\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\Gamma\|} [\|\Gamma\| \|\alpha\| + \|\delta\|]. \quad (1.4.14)$$

Отсюда легко вытекает неравенство

$$\frac{\|z - \alpha\|}{\|\alpha\|} \leq \frac{P_A}{1 - P_A \frac{\|\Gamma\|}{\|A\|}} \left[\frac{\|\Gamma\|}{\|A\|} - \frac{\|\delta\|}{\|f\|} \right]; \quad (1.4.15)$$

здесь $P_A = \|A\| \cdot \|A^{-1}\|$ - число обусловленности матрицы A . Формула (1.4.15) хорошо известна (см., например, Гавурий [1]). Она дает оценку относительной погрешности решения через число обусловленности матрицы A и относительные погрешности искажения $\|\Gamma\|/\|A\|$ и $\|\delta\|/\|f\|$. Оценка (1.4.15) очевидным образом переносится на произвольные уравнения вида (2.4.1), в которых A - производный оператор, действующий из одного банахова пространства в другое, при условии, что оба оператора A и A^{-1} ограничены. То же верно и для оценки (1.4.14), с той важной оговоркой, что ограниченность оператора A здесь необязательна.

ГЛАВА II

Погрешность аппроксимации

На протяжении всей книги мы будем пользоваться след. кривыми обозначениями и определениями. Норму в $W_2^{(3)}(\Omega)$ будем обозначать через $\|\cdot\|_{(3),\Omega}$ или, если это не может вызвать недоуменья, через $\|\cdot\|_{(3)}$; с означением $\|\cdot\|_3$ или $\|\cdot\|_{3,\Omega}$ мы сохраним для нормы в $L_3(\Omega)$.

Степенью одночлена $x_1^{d_1} x_2^{d_2} \dots x_m^{d_m}$ в R_m назовём мультииндекс $d = (d_1, d_2, \dots, d_m)$; степенью многочлена в R_m назовём верхнюю грань степеней соотвляющих его одночленов. Мультииндекс, все составляющие которого равны одному и тому же числу d , будем обозначать через \underline{d} . Таким образом, если $P(t)$ - многочлен в R_m , и по каждой из координат его степень не превосходит некоторого числа n , то степень $P(t)$ в R_m не превосходит мультииндекса \underline{n} .

Буквой C без индексов будем обозначать постоянные, вообще говоря, различные, точные значения некоторых для нас несущественны.

§ I. Классический метод Рунца. Оценки в энергетической норме.

1°. Напомним основы метода Рунца для линейных задач; подробно этот метод изложен в книге автора [3]. Пусть симметричный билинейный функционал $A(x, y)$ и соответствующий ему одноподный квадратичный функционал $A(x) = A(x, x)$ определены на некотором линейном множестве M . Допустим, что функционал $A(x)$ положителен: это значит, что $A(x) > 0$, если x отличен от нулевого элемента множества M . Превратим M в предгильбертово пространство, введя скалярное произведение $[x, y] = A(x, y)$ и норму $\|x\| = \sqrt{[x, x]} = \sqrt{A(x)}$. Замкнув M в этой норме, получим гильбертово пространство, которое обозначим через H_A . Для упрощения последующих записей примем, что пространство H_A - вещественное. В приложениях чаще встречается случай, когда M есть линейное плотное множество некоторого гильбертова /реже банахова/ пространства H , и $A(x, y) = (Ax, y)$, где A положительный /и симметричный - последняя оговорка необходима, если H - вещественное пространство/ оператор в H . В этом случае мы называем H_A энергетическим пространством оператора A , а скалярное произведение и норму в

H_A - соответственно энергетическим произведением и нормой.
 Пусть $f(x)$ - линейный функционал, ограниченный в H_A и
 определенный на всем этом пространстве. Поставим задачу о мини-
 муме функционала

$$F(x) = A(x) - 2f(x). \quad (\text{П.1.1.})$$

Эта задача теоретически решается элементарно. По известной тео-
 реме Ф.Рунда, в H_A существует один и только один элемент x_* ,
 такой, что $f(x) = [x, x_*], \forall x \in H_A$. Теперь $F(x) =$
 $= |x|^2 - 2[x, x_*] = |x - x_*|^2 - |x_*|^2$, и ясно, что функционал $F(x)$ дости-
 гает минимума при $x = x_*$.

2°. Приближенно можно решать задачу (П.1.1.), например, по
 классическому методу Рунда. Выбираем в H_A какое-нибудь конеч-
 номерное подпространство $H_A^{(n)}$ и ищем минимум функционала $F(x)$,
 или, что равносильно, минимум нормы $|x - x_*|$, на этом подпро-
 странстве. Элемент $x_*^{(n)} \in H_A^{(n)}$, на котором последний минимум дости-
 гается, существует и единственен. Очевидно, $x_*^{(n)}$ есть ортого-
 нальная проекция в H_A элемента x_* - точного решения задачи
 (П.1.1) - на подпространство $H_A^{(n)}$. Сравнивая это с общей схе-
 мой § I гл. I, видим, что в данном случае $X = H_A, X_n = H_A^{(n)}, P_n$
 есть тождественный оператор, соотношения \mathcal{U} и \mathcal{U}_n определяются
 формулами

$$x_*^{(n)} f^{(n)} = |x_*^{(n)}|^2 - 2f^{(n)}(x_*^{(n)}) = \min, x \in H_A^{(n)};$$

$$x_*^{(n)} f^{(n)} = |x_*^{(n)}|^2 - 2f^{(n)}(x_*^{(n)}) = \min, x^{(n)} \in H_A^{(n)}$$

и, наконец, $f^{(n)}$ есть сужение f на подпространство $H_A^{(n)}$.
 Отсюда ясно, что погрешность аппроксимации приближенного реше-
 ния по Рунду для задачи (П.1.1), равная $\rho_n = |x_* - x_*^{(n)}| = \min_{x \in H_A^{(n)}} |x_* - x|$
 совпадает с наилучшим приближением (в энергетической норме) точ-
 ного решения x_* элементами из пространства $H_A^{(n)}$. Величина ρ_n
 зависит от характера энергетической нормы и, в ряде важных случаев,
 от густоты элемента x_* ; при последовании последней полезно иметь в виду,
 что x_* есть решение уравнения Эйлера $\text{grad} F = 0$.

В ближайших параграфах настоящей главы мы применим высказан-
 ные в общих соображениях к крайним задачам для дифференциаль-
 ных уравнений.

Разумеется, применять метод Рунда целесообразно только если
 есть уверенность, что $|x_*^{(n)} - x_*| \xrightarrow{n \rightarrow \infty} 0$. Это последнее отно-

жение будет, очевидно, выполнено, если последовательность подпространств $\{H_A^{(n)}\}$ полна в H_A , иначе говоря, если по любому $\varepsilon > 0$ можно найти такое n_0 , что при любом $n \geq n_0$ существует элемент $x^{(n)} \in H_A^{(n)}$, удовлетворяющий неравенству $|x_n - x^{(n)}| < \varepsilon$.

В основополагающей работе Ритца [1] подпространства $H_A^{(n)}$ строятся так, что при возрастании n они расширяются. Более определенно, Ритц выбирает последовательность элементов $\{y_n\} \in H_A$, обладающую полнотой в H_A ; предполагается еще, что эти элементы, взятые в любом конечном числе, линейно независимы. За $H_A^{(n)}$ Ритц принимал подпространство с базисом (y_1, y_2, \dots, y_n) . Более поздние авторы называли систему $\{y_n\}$ координатной. Отметим, что Ритц не вводил энергетического пространства и требовал сходимости приближенных решений к точному в $C^{(3)}$ при подходящем выборе δ . Это ограничивало применение его метода и усложняло рассуждения. Энергетическая метрика была введена К.Фридрихсом; она была применена автором этой книги для наследования методов Ритца и Бубнова - Галеркина. Подробно об этом см. книгу автора [3].

Курант указал [1], что необязательно брать расширяющиеся подпространства $H_A^{(n)}$, и предложил также видоизменение метода Ритца: подпространства $H_A^{(n)}$ могут быть произвольные конечномерные подпространства энергетического пространства, подчиненные единственному требованию полноты их последовательности в H_A . Если выбрать в $H_A^{(n)}$ базис

$$(y_{n1}, y_{n2}, \dots, y_{nr}), \quad (П.1.2)$$

где $r = r(n)$ есть размерность $H_A^{(n)}$, то элементы y_{nk} обязательно принадлежат подпространству $H_A^{(n')}$, $n' > n$; таким образом, может случиться, что $H_A^{(n)}$ и $H_A^{(n')}$ имеют общим только нулевой элемент.

В случае расширяющихся подпространств будем говорить о классическом методе Ритца, в общем случае - об обобщенном методе Ритца. Классический метод Ритца имеет то преимущество, что погрешность $|x_n - x^{(n)}|$ убывает с ростом n (см. книгу автора [3]); следует, однако, напомнить, что обобщенный метод Ритца является идеальной основой метода конечных элементов - одного из самых распространенных в настоящее время методов решения краевых задач.

В заключение напомним, что если в $H_A^{(n)}$ выбран базис (П.1.2), то приближенное решение $x_n^{(n)}$ имеет вид

$$x_n^{(n)} = \sum_{k=1}^r a_k^{(n)} y_{nk} \quad (П.1.3)$$

причем коэффициенты $a_k^{(n)}$ определяются из алгебраической системы

$$\sum_{k=1}^N [y_{nk}, y_{nj}] a_k^{(n)} = f(y_{nj}), \quad j=1, 2, \dots, N. \quad (\text{П.1.4})$$

Последовательность базисов

$$(y_{n1}, y_{n2}, \dots, y_{nN}), \quad n=1, 2, \dots \quad (\text{П.1.5})$$

будем называть координатной системой обобщенного метода Рунге, а элементы y_{nk} - координатными элементами этого метода.

§ 2. Некоторые обобщения теоремы Джексона на функции многих переменных.

1°. Наиболее важное, как нам кажется, обобщение теоремы Джексона на функции многих переменных получила Харрик [1 - 3]. В ее работах [1, 2] рассмотрены более частные случаи; общая теорема доказана в работе [3]. Здесь мы формулируем основной результат этой последней статьи.

Пусть Ω - ограниченная область в R_m , τ и $\delta < \tau$ - натуральные числа. Далее, пусть $\varphi(t)$ - функция класса $C^{(\tau)}$ в некоторой области $\Omega_1 \supset \Omega$, причем ν достаточно велико, Ω_1 гомеоморфна шару, $\varphi|_{\partial\Omega} = 0$, $\text{grad } \varphi|_{\partial\Omega} \neq 0$ и $\varphi(t) \neq 0$, $t \notin \partial\Omega$.

Теорема П.2.1 (Харрик [3]). Пусть $x \in C^{(\nu)}(\bar{\Omega})$ и $\mathcal{D}^\alpha x|_{\partial\Omega} = 0$, $0 \leq |\alpha| \leq \delta - 1$. Тогда при любом натуральном n существует полином $P_n(t)$ степени $\leq n$ в R_m , удовлетворяющий неравенству

$$\|x - \varphi^\delta P_n\|_{C^{(\nu)}(\bar{\Omega})} \leq C \omega^{(\nu)}(x, \frac{1}{n}) n^{-(\nu - \bar{\nu})}; \quad 0 \leq \bar{\nu} \leq \tau. \quad (\text{П.2.1})$$

Здесь $\omega^{(\nu)}$ - наибольший из модулей непрерывности производных порядка ν от функции $x(t)$ в метрике $C(\bar{\Omega})$, а величина C не зависит ни от x , ни от n .

Доказательство теоремы П.2.1 основано на построении и исследовании операторов, аналогичных операторам Джексона. Заметим, что полагая в неравенстве (П.2.1) можно оценить C .

2°. Теорема П.2.2. (Шлошников [2]). Пусть область Ω и функция $\varphi(t)$ удовлетворяют условиям п.1°, $x \in W_2^{(\nu)}(\Omega)$ и $\mathcal{D}^\alpha x|_{\partial\Omega} = 0$, $0 \leq |\alpha| \leq \delta - 1$. Тогда для любого натурального n существует полином $P_n(t)$ степени $\leq n$ в R_m , удовлетворяющий неравенству

$$\|x - \varphi^\delta P_n\|_{L_2(\bar{\Omega})} \leq C \omega_{L_2}^{(\nu)}(x, \frac{1}{n}) n^{-(\nu - \bar{\nu})}, \quad 0 \leq \bar{\nu} \leq \tau, \quad (\text{П.2.2})$$

где C не зависит от x и n , а $\omega_{L_2}^{(\nu)}$ есть наибольший из L_2 -модулей непрерывности производных $\mathcal{D}^\alpha x$, $|\alpha| = \nu$.

Доказательство теоремы П.2.2 основано на том, что основные

свойства операторов, введенных Харрик, сохраняются при замене $C^{(k)}(\bar{\Omega})$ на $W_2^{(k)}(\Omega)$. Постоянная C в неравенстве (П.2.2) также можно оценить.

5°. Следствие П.2.1. Пусть граница $\partial\Omega$ представляет собой $(m-1)$ -мерную липшецеву поверхность, и $x \in W_2^{(\tau)}(\Omega)$. Тогда существует полином $Q_n(t)$ степени $\leq n$ в R_m , удовлетворяющий неравенству

$$\|x - Q_n\|_{(\bar{\tau}, \Omega)} \leq C_1 \max_{|d|=\tau} \|D^d x\|_{2, \Omega} n^{-(\tau-\bar{\tau})}, \quad 0 \leq \bar{\tau} \leq \tau - 1. \quad (\text{П.2.3})$$

Если поверхность $\partial\Omega$ достаточно гладкая, так что существует функция $\varphi(t)$ со свойствами, описанными в п.1°, то достаточно положить в (П.2.2) $\delta = 0$, заметив при этом, что $\omega_{L_2}^{(\tau)}(x, \delta) \leq 2 \max_{|d|=\tau} \|x^{(d)}\|_{2, \Omega}$. В общем случае продолжим функцию $x(t)$ с сохранением класса на все R_m так, чтобы она обращалась в нуль вне некоторого шара $\Omega_0 = \{t : |t| \leq R\}$ - такое продолжение возможно (см. Кальдерон [1]). Пусть $x^*(t)$ - продолженная функция; тем же символом $x^*(t)$ обозначим ее сужение на Ω_0 . Граница $\partial\Omega_0$ - сфера - есть бесконечно гладкая поверхность, и для функции $x^*(t)$, по оказанному выше, неравенство (П.2.3) справедливо: существует такой полином $Q_n(t)$ степени $\leq n$ в R_m , что

$$\|x^* - Q_n\|_{(\bar{\tau}, \Omega_0)} \leq C \max_{|d|=\tau} \|D^d x^*\|_{2, \Omega_0} n^{-(\tau-\bar{\tau})}. \quad (\text{П.2.4})$$

Функция $x^*(t)$ продолжает функцию $x(t)$ с сохранением класса, поэтому существует такая постоянная \mathfrak{E} ("постоянная продолжения"), что

$$|d| \leq \tau, \quad \|D^d x^*\|_{2, \Omega_0} \leq \mathfrak{E} \|D^d x\|_{2, \Omega};$$

кроме того, очевидно,

$$\|x^* - Q_n\|_{(\bar{\tau}, \Omega_0)} \geq \|x - Q_n\|_{(\bar{\tau}, \Omega)}.$$

Последние два соотношения позволяют заменить неравенство (П.2.4) неравенством (П.2.3) с постоянной $C_1 = \mathfrak{E}C$.

4°. Возвращаясь к теореме П.2.2. Пусть $m=1$ и область Ω есть промежуток $(0, 1)$. В этом случае можно положить $\varphi(t) = t(1-t)$ и неравенство (П.2.2) дает следующее: если $x \in W_2^{(\tau)}(0, 1) \cap \dot{W}_2^{(\tau)}(0, 1)$ то при любом натуральном n существует такой полином $P_n(t)$ степени $\leq n$, что

$$\|x - P_n(t)\|_{(\bar{\tau}, (0, 1))} \leq C \omega_{L_2}^{(\tau)}\left(x, \frac{1}{n}\right) n^{-(\tau-\bar{\tau})};$$

$$0 \leq \bar{v} \leq v, P_{n,v}(t) = t^3(1-t)^3 P_n(t). \quad (\text{П.2.5})$$

Точно так же из следствия П.2.1 вытекает, что если $x \in W_2^{(v)}(0,1)$, то при любом натуральном n существует такой полином $Q_n(t)$ степени $\leq n$, что справедливо неравенство

$$\|x - Q_n\|_{(\bar{v},(0,1))} \leq C \|D^v x\|_{2,(0,1)} \bar{v}^{-(v-\bar{v})}, \quad 0 \leq \bar{v} \leq v-1. \quad (\text{П.2.6})$$

§ 3. Обыкновенные дифференциальные уравнения.

1°. Рассмотрим краевую задачу для уравнения второго порядка

$$A x := -\frac{d}{dt} \left(p(t) \frac{dx}{dt} \right) + q(t)x = f(t); \quad 0 < t < 1, \quad (\text{П.3.1})$$

$$x(0) = x(1) = 0.$$

Примем следующие допущения: $p \in C^{(v)}[0,1]$, $q \in L_\infty(0,1)$,

$$f \in L_2(0,1); \quad p_0 \leq p(t) \leq p_1, \quad 0 < p_0 \leq p_1 < \infty, \quad q(t) \geq 0.$$

При этих допущениях оператор A — положительно определенный в $L_2(0,1)$; слабое решение задачи (П.3.1) совпадает с решением вариационной задачи

$$|x|^2 - 2 \int_0^1 f(t)x(t) dt = \min; \quad x(0) = x(1) = 0, \quad (\text{П.3.2})$$

где

$$|x|^2 = \int_0^1 \left\{ p(t) \left(\frac{dx}{dt} \right)^2 + q(t)x^2(t) \right\} dt. \quad (\text{П.3.3})$$

Пространство H_A состоит из функций $x(t)$, для которых интеграл (П.3.3) конечен и которые обращаются в нуль при $t=0$ и $t=1$.

Ниже нам будет полезна следующая оценка энергетической нормы:

если $q(t) \leq q_1$, то

$$|x|^2 \leq (p_1 + q_1/\pi^2) \|x'\|^2 =: C_0 \|x'\|^2, \quad (\text{П.3.4})$$

где $\|\cdot\|$ означает норму в $L_2(0,1)$. Доказательство очевидно следует из того, что наименьшее собственное число оператора $-d^2/dt^2$ при краевых условиях $x(0) = x(1) = 0$ равно π^2 .

2°. Са $H_A^{(n)}$ примем n -мерное подпространство с базисом $\{\sin k\pi t\}$, $1 \leq k \leq n$. Если $x^{(n)}$ произвольный элемент из $H_A^{(n)}$, то

$$\rho_n^2 = \|x_* - x_*^{(n)}\|^2 \leq \|x_* - x_*^{(n)}\|^2 \leq C_0 \|x_*' - x_*^{(n)'}\|^2. \quad (\text{П.3.5})$$

допустим, что $x_* \in W_2^{(v)}(0,1)$, $v \geq 3$. Для этого достаточно, чтобы $p \in W_2^{(v-1)}(0,1)$, $q, f \in W_2^{(v-2)}(0,1)$.

Функция $x_*(t)$ допускает разложение в ряд Фурье

$$x_*(t) = \sum_{k=1}^{\infty} \alpha_k \sin k\pi t.$$

За $x_*(t)$ примем отрезок этого ряда:

$$x^{(n)} = \sum_{k=1}^n a_k \sin k\pi t.$$

Оценим коэффициенты

$$a_k = 2 \int_0^1 x_*(t) \sin k\pi t dt.$$

Интегрируя по частям и принимая во внимание краевые условия (П.2.1), получим

$$a_k = - \frac{2}{k^3 \pi^3} \int_0^1 x_*'''(t) \cos k\pi t dt + \frac{2}{k^3 \pi^3} x_*''(t) \cos k\pi t \Big|_0^1.$$

Отсюда

$$|a_k| \leq \frac{C}{k^3} [\|x_*'''\|_{L_2} + \|x_*''\|_{L_\infty}].$$

Дальнейшее интегрирование по частям (возможное, если $\nu > 3$), не улучшает последней оценки. Теперь

$$\begin{aligned} \|x_*' - x^{(n)'}\|^2 &= \frac{1}{2} \sum_{k=n+1}^{\infty} k^2 \pi^2 a_k^2 \leq C [\|x_*'''\|_{L_2}^2 + \|x_*''\|_{L_\infty}^2] \times \\ &\times \sum_{k=n+1}^{\infty} \frac{1}{k^4} < C [\|x_*'''\|_{L_2}^2 + \|x_*''\|_{L_\infty}^2] \int_n^{\infty} \frac{dk}{k^4} = \\ &= C [\|x_*'''\|_{L_2}^2 + \|x_*''\|_{L_\infty}^2] \frac{1}{n^3} \end{aligned}$$

и формула (П.3.5) приводит к следующей оценке погрешности аппроксимации:

$$\rho_n \leq C [\|x_*'''\|_{L_2} + \|x_*''\|_{L_\infty}] n^{-3/2}. \quad (\text{П.3.6})$$

3°. По-прежнему пусть $x_* \in W_2^{(\nu)}(0,1)$. При $\nu > 3$ можно получить лучшую оценку аппроксимации, если за $H_A^{(n)}$ принять подпространство полиномов вида

$$Q_n(t) = t(1-t) \sum_{j=0}^k a_j t^j; \quad 0 \leq k \leq n.$$

В силу неравенства (П.2.5) существует такой полином $Q_n(t)$, что

$$\|x_* - Q_n\|_{L_2(\bar{\nu})} \leq C \omega_{L_2}^{(\nu)}(x_*, \frac{1}{n}) n^{-(\nu-\bar{\nu})}, \quad 0 \leq \bar{\nu} \leq \nu,$$

поэтому, если $x_*^{(n)}$ - приближенное решение из подпространства $H_A^{(n)}$, то

$$\rho_n = \|x_* - x_*^{(n)}\| \leq C \omega_{L_2}^{(\nu)}(x_*, \frac{1}{n}) n^{-\nu+1}. \quad (\text{П.3.7})$$

4°. Погрешность аппроксимации для обыкновенных дифференциальных уравнений в равномерных метриках изучали Крылов [1], автор [2], Бертрам [1], Берн-Буган [1], Леман [1] и Джикарани [2,3] в соболевских пространствах ту же погрешность исследовал Дичи [2] и Джикарани [3]. Вейлько [2] рассмотрел также задачу с

спектра.

§ 4. Уравнения эллиптического типа. Оценка в энергетической норме.

1°. Рассмотрим первую краевую задачу для эллиптического уравнения (или эллиптической системы уравнений) вида

$$\sum_{|\alpha|, |\beta|=0}^s (-1)^{|\alpha|} \mathcal{D}^\alpha (A_{\alpha\beta} \mathcal{D}^\beta x) = f(t), \quad x \in \Omega; \quad (П.4.1)$$

$$\mathcal{D}^\gamma x|_{\partial\Omega} = 0, \quad 0 \leq |\gamma| \leq s-1. \quad (П.4.2)$$

Пусть $\partial\Omega$ удовлетворяет требованиям п.1° § 2: поверхность $\partial\Omega$ можно задать уравнением $\varphi(t) = 0$, причем $\varphi \in C^{(s)}(\Omega_1)$, $\Omega_1 \supset \Omega$ и Ω_1 гомеоморфна шару; $\varphi(t) \neq 0$, $t \in \partial\Omega$ и $\text{grad } \varphi|_{\partial\Omega} \neq 0$.

Допустим, что оператор задачи (П.4.1) - (П.4.2) - симметричный и положительно определенный в $L_2(\Omega)$. Энергетическая норма для этого оператора определяется формулой

$$\|x\|^2 = \int_{\Omega} \sum_{|\alpha|, |\beta|=0}^s A_{\alpha\beta} \mathcal{D}^\alpha x \mathcal{D}^\beta x dt. \quad (П.4.3)$$

Будем считать, что коэффициенты $A_{\alpha\beta}$ суть функции от t , ограниченные и измеримые в $\bar{\Omega}$. Тогда, очевидно,

$$\|x\| \leq C_0 \|x\|_{(s), \Omega}. \quad (П.4.4)$$

Задачу (П.4.1) - (П.4.2) будем решать приближенно по методу Рунге, приняв за $H_A^{(n)}$ подпространство функций вида $\varphi^s(t) P_n(t)$, где $P_n(t)$ - полином степени $\leq n$ в R_m . В общем случае x_* - решение нашей задачи - принадлежит пространству $W_2^{(s)}(\Omega)$; допустим, что на самом деле оно более гладкое: пусть $x_* \in W_2^{(s)}(\Omega)$, где s достаточно большое число. Тогда существует такое произведение $\varphi^s(t) \tilde{P}_n(t) \in H_A^{(n)}$, что

$$\|x_* - \varphi^s \tilde{P}_n\| \leq C \omega_{L_2}^{(s)}(x_*, \frac{1}{n}) n^{-(s-1)}$$

Если $x_*^{(n)}$ - приближенное решение задачи (П.4.1) - (П.4.2) по Рунге, то, как это вытекает из сказанного в § 1, $\|x_* - x_*^{(n)}\| \leq \|x_* - \varphi^s \tilde{P}_n\| \leq C_0 \|x_* - \varphi^s \tilde{P}_n\|_{(s)}$ и, следовательно, $(s-1)$

$$\rho_n = \|x_* - x_*^{(n)}\| \leq C \omega_{L_2}^{(s)}(x_*, \frac{1}{n}) n^{-(s-1)} \quad (П.4.5)$$

2°. Пусть теперь дифференциальное уравнение (П.4.1) решается при естественных краевых условиях. По-прежнему будем считать, что оператор рассматриваемой задачи симметричный и положительно определенный в $L_2(\Omega)$, энергетическая норма удовлетворяет неравенству (П.4.4) и решение, которое мы по-прежнему обозначим через x_* , принадлежит к классу $W_2^{(s)}(\Omega)$. За $H_A^{(n)}$ примем

подпространство полиномов степени $\leq n$ в R_m . Используя следствие П.2.1 и повторяя рассуждения п.1⁰ данного параграфа, мы приходим к оценке

$$\|x_* - x_*^{(n)}\| \leq C \|x_*\|_{(x)} n^{-(\alpha-1)} \quad (\text{П.4.6})$$

Ту же оценку (П.4.6) можно получить, если краевые условия имеют вид (П.4.2) при $0 \leq |\alpha| \leq k-1$, $k < 5$, а остальные условия - естественные. В этом случае за $H_A^{(n)}$ следует принять подпространство функций $y^k(t) P_n(t)$, где $P_n(t)$ полином степени $\leq n$ в R_m .

3⁰. Условия, при которых точное решение $x_* \in W_2^{(\alpha)}(\Omega)$, см. Аймон, Дуглис, Ниренберг [1,2].

§ 5. Некоторые оценки погрешности производных. Литературные указания.

1⁰. Будем рассматривать задачу (П.4.1) - (П.4.2). В § 4 мы вывели оценку погрешности аппроксимации для приближенного решения этой задачи в энергетической метрике. По существу это оценка для производных, порядок которых равен половине порядка уравнения (П.4.1). Представляет значительный интерес получение аналогичных оценок для производных других порядков. Этому вопросу посвящен ряд работ, в которых рассматриваются координатные функции, описанные в п.2 § 4.

Ильин [1,3] рассмотрел основные краевые задачи для невырожденного эллиптического уравнения второго порядка; для этих задач получены оценки нормы $\|x_* - x_*^{(n)}\|_{C(\bar{\Omega})}$ при условии, что известен порядок убывания некоторых величин A_n , определенным образом зависящих от приближенного решения $x_*^{(n)}$. В аналогичном предположении получены оценки погрешности аппроксимации в $C(\bar{\Omega})$ и в $C^{(k)}(\bar{\Omega})$ в случае первой краевой задачи для бигармонического уравнения. В [3] Ильин улучшил оценки для эллиптического уравнения второго порядка. Для того же уравнения Харрик [1] получала оценку величин A_n ; аналогичная оценка для бигармонического уравнения дана в работе Харрик [2]. В статье Ильина [5] рассмотрены эллиптические уравнения любого четного порядка 2β и получены оценки погрешности аппроксимации в $W_p^{(\beta)}(\Omega)$, $1 < p < \infty$, β - не очень большое четносрицательно целое число.

2⁰. В статье Канторовича [3] (см. также монографии Канторовича и Акилова [1]) предложена некоторая общая схема исследования погрешности аппроксимации проекционных методов в банаховых

пространствах. В статье автора [II] предложено^{I)} некоторое видоизменение этой схемы, которое позволило получить оценку погрешности аппроксимации метода Рунге в $W_2^{(\delta)}$, где $\delta > 2\delta$ - не очень большое натуральное число. Следует указать, что оценка этой последней статьи хуже оценок статьи Ильина [5], соответствующих значению $\rho = 2$. Другое видоизменение схемы Канторовича было использовано Шапошниковой [1], которая таким способом получила оценку погрешности аппроксимации для производных порядка $\delta < 2\delta$. Ее оценки хуже оценок Ильина [5] при $\rho = 2$ для $\delta > \delta$ и лучше для $\delta < \delta$. В статье Шапошниковой [3] получены оценки в некоторых банаховых пространствах; показано, в частности, что для $\delta < \delta$ оценки Ильина [5] можно улучшить и при $\rho \neq 2$.

3°. В статье автора [32] метод его работы [II] и работы Шапошниковой [1] перенесен на несамосопряженные эллиптические задачи, допускающие решение по методу Бубнова - Галеркина. Благодаря использованию уточненных марковских неравенств удалось улучшить результаты названных выше работ автора [II] и Шапошниковой [1]; новые оценки совпадают с оценками Ильина (при $\rho = 2$) для $\delta > \delta$ и с оценками Шапошниковой [1] при $\delta < \delta$.

Ниже, в § 6, будут изложены необходимые марковские неравенства, в §§ 7, 8 - результаты работ автора [II, 32] и Шапошниковой [1] для самосопряженных эллиптических задач; в § 9 будут рассмотрены несамосопряженные эллиптические задачи.

§ 6. Некоторые марковские неравенства

1°. Пусть $P_n(t)$ - полином степени $\leq n$ от вещественной переменной t на промежутке $a \leq t \leq b$. Обозначим через $q_k(t)$ $k = 0, 1, 2, \dots$, полиномы степени k , ортонормированные в $L_2(a, b)$ эти полиномы отличаются лишь постоянным множителем $[2(b-a)]^{k/2}$ от нормированных полиномов Лемандра, преобразованных к промежутку (a, b) . Справедливо тождество

$$P_n(t) = \sum_{k=0}^n a_k q_k(t); \quad a_k = \text{const.}$$

Оценим

$$P_n^2(t) \leq \sum_{k=0}^n a_k^2 \sum_{k=0}^n q_k^2(t) = \|P_n\|_{L_2(a,b)}^2 \sum_{k=0}^n q_k^2(t).$$

^{I)} Результаты статьи автора [II], вышедшей из печати только в 1970 г., были доложены в 1968 г. на симпозиуме в честь 50-летия акад. С. Л. Соболева.

Далее, $q_k(t) = [2/(b-a)]^{1/2} p_n(\tau)$, где $p_n(\tau)$ - полином Лежандра, ортонормированные на промежутке $(-1, +1)$, и $\tau = (b-a)^{-1}(2t-a-b)$. Так как $\max_{|\tau| \leq 1} p_n(\tau) = \sqrt{(2n+1)/2}$, то

$$\max_{t \in [a,b]} \sum_{k=0}^n q_k^2(t) = \frac{2}{b-a} \sum_{k=0}^n \frac{2n+1}{2}$$

и, следовательно,

$$\|P_n\|_{C[a,b]} \leq Cn \|P_n\|_{L_2(a,b)}. \quad (П.6.1)$$

Это неравенство хорошо известно.

2°. Выберем марковское неравенство для полиномов от многих переменных в L_2 . В статье Хилле, Стге и Тамаркина [1] (см. также Бари [1]) для полиномов от одной вещественной переменной установлено следующее неравенство: если $P_n(t)$ - такой полином степени n , то

$$\|P_n'\|_{L_p(a,b)} \leq Cn^2 \|P_n\|_{L_p(a,b)}, \quad (П.6.2)$$

где C зависит только от p и от $b-a$. Мы будем пользоваться этим неравенством для значения $p = 2$.

Неравенство (П.6.2) очевидным образом обобщается на полиномы степени $\leq n$ в R_m , если отрезок (a, b) заменяется параллелепипедом. Действительно, можно считать, что ребра параллелепипеда параллельны координатным осям. Пусть $d = (d_1, d_2, \dots, d_m)$ - мультииндекс, причем $d_k = \delta_{jk}$, где индекс j фиксирован, и пусть для точек j -го ребра параллелепипеда $a_j \leq t_j \leq b_j$. По неравенству (П.6.2)

$$\int_{a_j}^{b_j} |D^d P_n(t)|^2 dt_j \leq Cn^4 \int_{a_j}^{b_j} |P_n(t)|^2 dt.$$

Интегрируем по остальным координатам, получим (Π - параллелепипед)

$$|d| = 1, \|D^d P_n\|_{L_2(\Pi)} \leq Cn^2 \|P_n\|_{L(\Pi)}. \quad (П.6.3)$$

Отсюда

$$\forall d, \|D^d P_n\|_{L_2(\Pi)} \leq Cn^{2|d|} \|P_n\|_{L(\Pi)}. \quad (П.6.4)$$

3°. Рассмотрим теперь в R_m некоторую конечную область Ω с границей, удовлетворяющей условиям п.1° § 2. Представим Ω в виде $\Omega = \Omega' \cup \Omega''$; $\Omega' \cap \Omega'' = \emptyset$. Здесь Ω'' - внутренняя полоса ширины δ и δ - достаточно малое фиксированное число. Каждую точку замкнутой области Ω сделаем центром куба, ребра которого параллельны осям координат и который целиком вместе со своей границей лежит в Ω . По лемме 10-го абзаца выберем конечное

число таких кубов, образующих покрытие для Ω' ; пусть это будут кубы Q_1, Q_2, \dots, Q_k . В силу неравенства (П.6.4)

$$\forall d, \|D^{\alpha} P_n\|_{2, Q_j} \leq C n^{2|\alpha|} \|P_n\|_{2, Q_j}.$$

Прооуммируем это по j . Слева заменим полученную оумму меньшей величиной $\|D^{\alpha} P_n\|_{2, \Omega'}$, а справа - большей величиной $k \|P_n\|_{2, \Omega}$. В результате получим

$$\forall d, \|D^{\alpha} P_n\|_{2, \Omega'} \leq C n^{2|\alpha|} \|P_n\|_{2, \Omega}. \quad (\text{П.6.5})$$

4°. Теперь рассмотрим пограничную полосу Ω_{δ} . Ширину δ этой полосы выберем следующим образом. Пусть $t_0 \in \partial\Omega$. Построим сферу с центром в t_0 , обладающую тем свойством, что любая прямая, параллельная нормали к $\partial\Omega$ в точке t_0 , пересекает часть Γ_0 поверхности $\partial\Omega$, заключенную внутри упомянутой сферы, не более одного раза. Известно, (см., например, книгу автора [10]), что радиус d этой сферы можно выбрать независимым от t_0 . Участью Γ_0 границы $\partial\Omega$ можно задать уравнением вида $\xi_m = q(\xi)$, где $\xi = (\xi_1, \xi_2, \dots, \xi_{m-1})$, ось ξ_m направлена по нормали к $\partial\Omega$ в точке t_0 и начало новой системы координат совпадает с t_0 . Так как функция $\varphi(t)$ (см. п. I° § 2) достаточно гладкая, причем $\text{grad } \varphi|_{\partial\Omega} \neq 0$, то функция $q(\xi)$ также достаточно гладкая, и можно найти столь малое $\delta > 0$, чтобы лежащая в Ω_{δ} замкнутая область

$$E_{t_0, \delta} = \{t: -\delta \leq \xi_i \leq \delta, 1 \leq i \leq m-1;$$

$$q(\xi) - \delta \leq \xi_m \leq q(\xi)\}$$

находилась внутри сферы радиуса d с центром в t_0 .

Введем координаты

$$\eta_i = \xi_i, 1 \leq i \leq m-1, \eta_m = -\xi_m + q(\xi).$$

В координатах $\eta_1, \eta_2, \dots, \eta_m$ область $E_{t_0, \delta}$ переходит в параллелепипед

$$F_{\delta} = \{\eta: -\delta \leq \eta_i \leq \delta, 1 \leq i \leq m-1; 0 \leq \eta_m \leq \delta\}.$$

Очевидно, якобиан преобразования от координат ξ_i к координатам η_i по абсолютной величине равен единице, а производные первого порядка от x по t_i оцениваются через те же производные по η_j , $i, j = 1, 2, \dots, m$, и наоборот.

Каждой точке $t_0 \in \partial\Omega$ можно привести в соответствие область $E_{t_0, \delta}$, и эти области образуют покрытие для Ω_{δ} . По лемме Бореля можно выбрать конечное число областей $E_{t_0, \delta}$, также образу-

ных покрытие для Ω_δ . Для каждого из параллелепипедов F_δ справедливо неравенство (П.6.4), из которого в свою очередь следует новое неравенство

$$\forall d, \|D^d P_n\|_{2, \Omega_\delta} \leq C \|P_n\|_{2, \Omega} n^{2|d|}.$$

Сложив это с (П.6.5), получим искомое марковское неравенство

$$\forall d, \|D^d P_n\|_{2, \Omega} \leq C \|P_n\|_{2, \Omega} n^{2|d|}. \quad (\text{П.6.6})$$

Нетрудно получить аналогичное неравенство и для $p \neq 2$.

5°. Докажем, что справедливо неравенство

$$\|y^{\delta-1} P_n\|_{2, \Omega} \leq C n^2 \|y^\delta P_n\|_{2, \Omega}. \quad (\text{П.6.7})$$

Начнем со случая одного измерения. В этом случае можно свести дело к сравнению интегралов

$$\int_0^\delta P_n^2(t) dt, \quad \int_0^\delta t^2 P_n^2(t) dt,$$

где $\delta > 0$ — фиксированное число. Интегрируя по частям, получим

$$\int_0^\delta P_n^2(t) dt = \delta P_n^2(\delta) - 2 \int_0^\delta t P_n(t) P_n'(t) dt.$$

Оценим правую часть. Прежде всего, по формуле (П.6.1),

$$\delta P_n^2(\delta) \leq \frac{1}{\delta} \|t P_n\|_{C[0, \delta]}^2 \leq C n^2 \|t P_n\|_{2, (0, \delta)}^2.$$

далее,

$$\begin{aligned} |2 \int_0^\delta t P_n(t) P_n'(t) dt| &\leq \frac{1}{\varepsilon} \int_0^\delta t^2 P_n^2(t) dt + \varepsilon \int_0^\delta P_n'^2(t) dt \leq \\ &\leq \frac{1}{\varepsilon} \int_0^\delta t^2 P_n^2(t) dt + C \varepsilon n^4 \int_0^\delta P_n^2(t) dt. \end{aligned}$$

Выбрав $\varepsilon = (2Cn^4)^{-1}$, найдем искомую оценку:

$$\|P_n\|_{2, (0, \delta)} \leq C n^2 \|t P_n\|_{2, (0, \delta)}. \quad (\text{П.6.8})$$

Обратимся к общему случаю. Внутри Ω оценка тривиальна: если Ω' — внутренняя подобласть и $|y(t)| \geq \beta$, $\forall t \in \Omega'$, то

$$\|y^{k-1} P_n\|_{2, \Omega'} \leq \frac{1}{\beta} \|y^k P_n\|_{2, \Omega'}. \quad (\text{П.6.9})$$

Достаточно повсюду рассмотреть область вида $E_{t_0, \delta}$ и сравнить интегралы

$$\int_{E_{t_0, \delta}} y^{2k-2}(t) P_n^2(t) dt, \quad \int_{E_{t_0, \delta}} y^{2k} P_n^2(t) dt.$$

делаем замену переменных, перейдя от переменных t_i к переменным η_i (см. п.4°). Область $E_{t_0, \delta}$ перейдет в параллелепипед F_δ , для этого интеграл вида

$$\int_{E_{t_0, \delta}} y^{2k}(t) P_n^2(t) dt \quad (П.6.10)$$

оценится сверху и снизу через интеграл

$$\int_{J_\delta} \eta_m^{2k} P_n^2(t) d\eta \quad (П.6.11)$$

с постоянными, не зависящими от n . Выражение $\eta_m^{k-1} P_n(t)$ есть полином степени $\leq n+k-1$ от η_m , и по неравенству (П.3.6)

$$\int_0^\delta \eta_m^{2k-2} P_n^2(t) d\eta_m \leq Cn^4 \int_0^\delta \eta_m^{2k} P_n^2(t) d\eta_m.$$

Интегрируя это по η_i , $1 \leq i \leq m-1$, в пределах от $-\delta$ до $+\delta$ и пользуясь эквивалентностью интегралов (П.6.10) и (П.6.11), найдем, что

$$\int_{E_{t_0, \delta}} y^{2k-2}(t) P_n^2(t) dt \leq Cn^4 \int_{E_{t_0, \delta}} y^{2k}(t) P_n^2(t) dt.$$

Отсюда и из неравенства (П.3.9) вытекает оценка (П.3.7).

§ 7. Оценки аппроксимации старших производных

1°. Уже было упомянуто, что Л.В. Канторович предложил некоторую схему построения и последования проекционных методов. Мы изложим здесь эту схему в условиях не совсем общих, но достаточных для того, чтобы ее можно было применить к методу Рунца.

Рассмотрим уравнение

$$Ax = f, \quad (П.7.1)$$

где A — линейный оператор, взаимно однозначно отображающий банахово пространство X на банахово пространство Y , тогда оба оператора A и A^{-1} ограничены. Построим полную в X последовательность конечномерных пространств X_n ; напомним, что это означает следующее:

$$\forall x \in X, \quad \varepsilon_n(x) = \inf_{x^{(n)} \in X_n} \|x - x^{(n)}\|_X \xrightarrow{n \rightarrow \infty} 0. \quad (П.7.2)$$

Положим $Y_n = AX_n$; для каждого n выберем оператор Q_n , проектирующий Y на Y_n . Приближенное решение уравнения (П.7.1) будем строить как элемент $x^{(n)} \in X_n$, удовлетворяющий уравнению

$$Q_n Ax^{(n)} = Q_n f. \quad (П.7.3)$$

Если x_* и $x_*^{(n)}$ - точные решения задач (П.7.1) и (П.7.3) соответственно, то, как доказано в работах Канторовича, цитированных выше,

$$\|x_* - x_*^{(n)}\|_X \leq (1 + \|A^{-1} Q_n A\|_X) \delta_n(x_*). \quad (\text{П.7.4})$$

Следовательно оценка

$$\|x_* - x_*^{(n)}\|_X \leq C \|Q_n\|_Y \delta_n(x_*). \quad (\text{П.7.5})$$

2°. В статье [11] автор предложил некоторое видоизменение схемы Канторовича. Пусть $X \subset Y \subset Z$; здесь Y и Z - гильбертовы пространства, X - банахово пространство; знак \subset означает плотное вложение. В уравнении (П.7.1) будем рассматривать A также как оператор в Z и будем считать, что он положительно определен, так что

$$\forall x \in \mathcal{D}(A), (Ax, x)_Z \geq r^2 \|x\|_Z^2; \quad r^2 = \text{const} > 0;$$

в то же время, в соответствии со сказанным в п.1°, будем считать, что A взаимно однозначно отображает X на Y .

Имеем полную в X последовательность подпространств $\{X_n\}$; пусть $\{y_{n1}, y_{n2}, \dots, y_{nk}\}$ - базис в X_n . Положим $Y_n = AX_n$ и определим оператор Q_n , отображающий Y на Y_n , соотношением

$$\forall f \in Y, (f - Q_n f, y_{nj})_Z = 0, \quad j = 1, 2, \dots, k \quad (\text{П.7.6})$$

Докажем, что Q_n есть проектор из Y на Y_n . По определению, $Q_n: Y \rightarrow Y_n$, поэтому, если $f \in Y$, то $Q_n f \in Y_n$ и, следовательно

$$Q_n f = \sum_{k=1}^k a_k^{(n)} A y_{nk}.$$

Далее, по тому же определению, $(Q_n f - f, y_{nj}) = 0$, или

$$\sum_{k=1}^k a_k^{(n)} (A y_{nk}, y_{nj})_Z = (f, y_{nj})_Z, \quad 1 \leq j \leq k. \quad (\text{П.7.7})$$

Оператор A положительно определен в Z , поэтому и матрица системы (П.7.7) положительно определенная; эта система разрешима единственным образом, и оператор Q_n определен для всех натуральных n . Остается доказать, что $Q_n^2 = Q_n$. Если $f \in Y$, то $Q_n f \in Y_n \subset Y$ и $Q_n^2 f \in Y_n$, или

$$Q_n^2 f = \sum_{k=1}^k b_k^{(n)} A y_{nk}.$$

Далее, $(Q_n^2 f - Q_n f, y_{nj}) = 0$, поэтому $(Q_n^2 f, y_{nj}) = (Q_n f, y_{nj}) = (f, y_{nj})$, или

$$\sum_{k=1}^N b_k^{(n)} (A y_{nk}, y_{nj}) = (f, y_{nj}), \quad j = 1, 2, \dots, N.$$

Система (П.7.7) разрешима единственным образом, поэтому $b_k^{(n)} = a_k^{(n)}$ и $Q_n^2 f = Q_n f$.

Заменяя f на $Ax^{(n)}$ под знаком проектора Q_n в (П.7.6), получим систему уравнений Бубнова - Галеркина для уравнения (П.7.1); так как A - положительно определенный оператор в Z , то (П.7.6), или равносильная система (П.7.7), есть также и система Рунца для того же уравнения. Поскольку $Q_n f \in Y_n$, то

$$Q_n f = \sum_{k=1}^N a_k^{(n)} A y_{nk}, \quad a_k^{(n)} = \text{const}$$

и, следовательно,

$$\|Q_n f\|_Y^2 = \left\| \sum_{k=1}^N a_k^{(n)} A y_{nk} \right\|_Y^2 \leq \Lambda_n^{(n)} \sum_{k=1}^N |a_k^{(n)}|^2,$$

где $\Lambda_n^{(n)}$ - наибольшее собственное число матрицы скалярных произведений $(A y_{nk}, A y_{nj})_Y$; $j, k = 1, 2, \dots, N$. Заметим, что $a_k^{(n)}$ суть коэффициенты Рунца приближенного решения $x_*^{(n)}$:

$$x_*^{(n)} = \sum_{k=1}^N a_k^{(n)} y_{nk}.$$

Обозначая, как обычно, через $[,]$ и $\|\cdot\|$ скалярное произведение и норму в энергетическом пространстве Z оператора A , получаем из последнего тождества

$$\begin{aligned} \|x_*^{(n)}\|^2 &= \sum_{j,k=1}^N a_j^{(n)} a_k^{(n)} [y_{nj}, y_{nk}] = \\ &= \sum_{j,k=1}^N a_j^{(n)} a_k^{(n)} (A y_{nj}, A y_{nk})_Z \geq \lambda_n^{(n)} \sum_{k=1}^N |a_k^{(n)}|^2; \end{aligned}$$

здесь $\lambda_n^{(n)}$ - наименьшее собственное число матрицы скалярных произведений $(A y_{nj}, A y_{nk})_Z$; $j, k = 1, 2, \dots, N$. Теперь (см. книгу автора [3])

$$\sum_{k=1}^N |a_k^{(n)}|^2 \leq \frac{\|x_*^{(n)}\|^2}{\lambda_n^{(n)}} \leq \frac{\|x_*\|^2}{\lambda_n^{(n)}} = \frac{\|f\|_Z^2}{r^2 \lambda_n^{(n)}}.$$

Из вложения $Y \subset Z$ вытекает существование такой постоянной C , что $\|f\|_Z = C \|f\|_Y$, и из (П.7.8) вытекает искомая оценка:

$$\|Q_n\|_{y \rightarrow y_n} \leq C \sqrt{\Lambda_n^{(n)} / \lambda_n^{(n)}}. \quad (П.7.9)$$

Если базис подпространства $H_A^{(n)}$ ортонормирован в Z_A , то $\lambda_n^{(n)} = 1$, и последняя оценка упрощается:

$$\|Q_n\|_{y \rightarrow y_n} \leq C \sqrt{\Lambda_n^{(n)}}. \quad (П.7.10)$$

Из формулы (П.7.5) получается оценка погрешности аппроксимации в X -норме:

$$\|x_* - x_*^{(n)}\|_X \leq C \sqrt{\Lambda_n^{(n)}} \varepsilon_n(x_*). \quad (П.7.11)$$

3°. Возвращаясь к задаче (П.4.1) - (П.4.2). Допустим, что система (П.4.1) равномерно эллиптическая в Ω , а оператор рассматриваемой задачи положительно определен в $L_2(\Omega)$. Решение $u_*(t)$ будет сколь угодно гладким, если достаточно гладкими будут функции $g(t)$, $f(t)$, $A_{\alpha\beta}(t)$. Будем считать, что $x_* \in W_2^{(l)}(\Omega)$, где l достаточно большое число.

Рассмотрим сначала более простой случай, когда $g(t)$ есть полином. Задачу (П.4.1) - (П.4.2) будем решать по Рунту, взяв в качестве координатных функций полиномы $g^s(t) P_n^{(j)}(t)$, где $P_n^{(j)}$ - полиномы степени $\leq n$ в R_m , тогда приближенное решение имеет вид $x_*^{(n)}(t) = g^s(t) P_n(t)$, где $P_n(t)$ - также полином степени $\leq n$ в R_m .

Введем в рассмотрение пространства $X = W_2^{(s)}(\Omega)$, $Y = W_2^{(l)}(\Omega)$, $Z = L_2(\Omega)$; $l \geq 0$, $s = 2s + l \leq r$. Если сверху здесь означает мощность пространства $W_2^{(s)}(\Omega)$, определяемое условиями (П.4.2).

Координатные функции $y_{nj}(t)$ будем считать ортонормированными в $W_2^{(s)}(\Omega)$ - тогда, как нетрудно проверить (подробнее см. книгу автора [8]) числа $\lambda_n^{(l)}$ положительно ограничены снизу и величина оценка (П.7.11).

Пусть K - число функций y_{nj} при фиксированном n . Положим $\tau = (\tau_1, \tau_2, \dots, \tau_K)$ и $\|\tau\|_2^2 = |\tau_1|^2 + |\tau_2|^2 + \dots + |\tau_K|^2$. Тогда

$$\lambda_n^{(l)} = \max_{\|\tau\|_2=1} \|A \varphi_n(\cdot, \tau)\|_Y^2, \quad \text{где } \varphi_n(t, \tau) = \sum_{j=1}^K \tau_j y_{nj}(t).$$

Для любой функции $u \in W_2^{(s)}(\Omega)$ справедливо неравенство

$\|Au\|_y \leq C \|u\|_x$; в частности, $\|A\varphi_n\|_y \leq C \|\varphi_n\|_x =$
 $= C \|\varphi_n\|_{(\tilde{\delta}), \Omega}$. Но $\varphi_n(t, \tau)$ есть полином степени $\leq n + \text{const}$ в R_m , и то неравенству (П.6.3)

$$\begin{aligned} \|\varphi_n\|_{(\tilde{\delta}), \Omega} &\leq C n^{2(\tilde{\delta}-1)} \|\varphi_n\|_{(\tilde{\delta}), \Omega} \leq \\ &\leq C n^{2(\tilde{\delta}-1)} \sum_{k=1}^n |\tau_k|^2 = C n^{2(\tilde{\delta}-1)}. \end{aligned}$$

Теперь

$$\Lambda_n^{(n)} \leq C n^{4(\tilde{\delta}-1)}. \quad (\text{П.7.12})$$

Остается оценить наилучшее приближение

$$\mathcal{E}_n(x_*) = \inf_{x \in \chi_n} \|x_* - x\|_X = \inf_{x^{(n)} \in \chi_n} \|x_* - x^{(n)}\|_{(\tilde{\delta}), \Omega}.$$

Функция $x^{(n)}$ есть полином степени $\leq n + \text{const}$ в R_m , а

$x_* \in W_2^{(\tilde{\delta})}(\Omega)$; по теореме П.2.2 существует такой полином $\bar{x}^{(n)} \in \chi_n$, что

$$\|x_* - \bar{x}^{(n)}\|_{(\tilde{\delta}), \Omega} \leq C n^{-(\tilde{\delta}-1)} \omega_{L_2}^{(\tilde{\delta})}(x_*; \frac{1}{n}).$$

Это значит, что $\mathcal{E}_n(x_*) \leq C n^{-(\tilde{\delta}-1)} \omega_{L_2}^{(\tilde{\delta})}(x_*; \frac{1}{n})$.

Теперь из формул (П.7.5) и (П.7.12) следует:

$$\rho_n = \|x_* - \mathcal{E}_n^{(n)}\|_{(\tilde{\delta}), \Omega} \leq C n^{-(\tilde{\delta}-1)2\tilde{\delta}} \omega_{L_2}^{(\tilde{\delta})}(x_*; \frac{1}{n}); \quad (\text{П.7.13})$$

это и есть оценка погрешности аппроксимации для производных порядка $\tilde{\delta} \geq 2$.

4°. Вернемся к обскому случаю и будем считать, что функция $(\cdot)^{(t)}$ достаточно гладкая, но не обязательно полином. Введем оценку величины $\Lambda_n^{(n)} = \max_{\tau \in T} \|A\varphi(\cdot, \tau)\|_y$. Для лю-

бой функции $u \in W_2^{(\tilde{\delta})}(\Omega)$ справедливо неравенство $\|Au\|_{(e)} \leq C \|u\|_{(\tilde{\delta})}$; в частности, $\|A\varphi_n\|_{(e)} \leq C \|\varphi_n\|_{(\tilde{\delta})}$. Формулу в $W_2^{(\tilde{\delta})}(\Omega)$ можно задать формулой

$W_2^{(\tilde{\delta})}(\Omega)$ можно задать формулой

$$\|u\|_{(3)} = \sum_{\sigma, \zeta} \|\mathcal{D}^{\sigma+\zeta} u\|_2;$$

суммирование производится по всем мультииндексам σ и ζ , таким, что $|\sigma| \geq \delta + \nu$ и $|\zeta| \leq \delta$. Докажем неравенство

$$\|\mathcal{D}^{\sigma+\zeta} \varphi_n\|_2 \leq C n^{2|\sigma|} \|\varphi_n\|_{(3)}. \quad (\text{П.7.14})$$

Имеем $\varphi_n(t, \tau) = y^\delta(t) P_n(t, \tau)$, где $P_n(t, \tau)$ - полином степени $\leq \underline{n}$ в R_m . Построим полином $q_n(t)$ степени $\leq \underline{n}$ в R_m , такой, что

$$\forall \nu, |\nu| \leq |\sigma + \zeta|, \mathcal{D}^\nu [y^\delta(t) - q_n(t)] = O(\epsilon n^{-\nu+\nu});$$

это возможно в силу теоремы П.2.1. Будем считать, что ν достаточно велико. Далее,

$$\begin{aligned} \|\mathcal{D}^{\sigma+\zeta} \varphi_n(\cdot, \tau)\|_2 &\leq \|\mathcal{D}^{\sigma+\zeta} (q_n P_n(\cdot, \tau))\|_2 + \\ &+ \|\mathcal{D}^{\sigma+\zeta} (y^\delta - q_n) P_n(\cdot, \tau)\|_2. \end{aligned} \quad (\text{П.7.15})$$

Выражение $\mathcal{D}^\zeta (q_n(t) P_n(t, \tau))$ есть полином степени $\leq 2\underline{n} - \zeta \leq 2\underline{n}$ в R_m , и по неравенству (П.6.4)

$$\|\mathcal{D}^{\sigma+\zeta} (q_n P_n(\cdot, \tau))\|_2 \leq C n^{2|\sigma|} \|\mathcal{D}^\zeta (q_n P_n(\cdot, \tau))\|_2. \quad (\text{П.7.16})$$

$$\begin{aligned} &\|\mathcal{D}^\zeta (q_n P_n(\cdot, \tau))\|_2 \leq \\ &\leq \|\mathcal{D}^\zeta (y^\delta P_n(\cdot, \tau))\|_2 + \|\mathcal{D}^\zeta [(y^\delta - q_n) P_n(\cdot, \tau)]\|_2. \end{aligned} \quad (\text{П.7.17})$$

Оценим второе слагаемое:

$$\begin{aligned} \|\mathcal{D}^\zeta [(y^\delta - q_n) P_n(\cdot, \tau)]\|_2 &= \left\| \sum_{d \leq \zeta} \mathcal{D}^d (y^\delta - q_n) \mathcal{D}^{\zeta-d} P_n(\cdot, \tau) \right\|_2 \leq \\ &\leq \frac{C}{n^{|\sigma-|\zeta|}} \sum_{d \leq \zeta} \|\mathcal{D}^{\zeta-d} P_n(\cdot, \tau)\|_2 \leq C n^{-\nu+3|\zeta|} \|P_n(\cdot, \tau)\|_2 \leq \end{aligned}$$

$$\leq C n^{-\sigma+3|k|+2\delta} \|y^\delta P_n(\cdot, \tau)\|_2 \leq C \|y^\delta P_n(\cdot, \tau)\|_2;$$

здесь используются неравенства (П.6.6) и (П.6.II), а также то, что U достаточно велико. Соболевские теоремы вложения дают неравенство

$$\|y^\delta P_n(\cdot, \tau)\|_2 \leq C \|y^\delta P_n(\cdot, \tau)\|_{(3)},$$

а, следовательно, и неравенство

$$\|D^z [(y^\delta - q_n) P_n(\cdot, \tau)]\|_2 \leq C \|y^\delta P_n(\cdot, \tau)\|_{(3)}. \quad (\text{П.7.18})$$

Точно так же находим

$$\|D^{z+\beta} [(y^\delta - q_n) P_n(\cdot, \tau)]\|_2 \leq C \|y^\delta P_n(\cdot, \tau)\|_{(3)}. \quad (\text{П.7.19})$$

Из неравенств (П.7.15-19) очевидно следует (П.7.14), из которого в свою очередь следует, что

$$\|\varphi_n\|_{(3)}^2 \leq C n^{4\tilde{\delta}-2\delta} \|\varphi_n\|_{(3)}^2 \leq C n^{4\tilde{\delta}-2\delta}.$$

В силу теоремы П.2.1, $\mathcal{E}_n(x_*) \leq C n^{-\tau+\delta} \omega_{L_2}^{(\nu)}(x_*, \frac{1}{n})$. Послед-

ние два неравенства вместе с формулой (П.7.10) приводят к оценке (П.7.13).

5°. Если уравнение (П.4.1) решать при естественных краевых условиях, то можно получить оценку, близкую к (П.7.13). За координатные функции можно взять полиномы, при этом все предшествующие рассуждения в существенном сохраняются, но для наилучшего приближения получается (П.2.3), и окончательно

$$\rho_n = \|x_* - x_*^{(n)}\|_{(3)} \leq C n^{-\tau+3\tilde{\delta}-2\delta}. \quad (\text{П.7.20})$$

§ 8. Оценка аппроксимации младших производных

1°. Такой рода оценки получила, в частности, Шапошникова [I], используя еще одно видоизменение схемы Канторовича. По-прежнему к пространствам X и Y добавляется еще одно пространство Z и по-прежнему пространства Y и Z гильбертовы, но, в отличие от § 7, рассматриваемые пространства связаны соотношением

$X \subset Z \subset Y$. Допустим, что существует такое сужение \hat{A} оператора A задачи (П.7.1), что $\mathcal{D}(\hat{A}) \subset Z$, $R(\hat{A}) \subset Z$ и \hat{A} - положительно определенный оператор в Z . Энергетическое пространство оператора \hat{A} обозначим через \hat{Z} . Возможны два случая взаимного расположения пространств:

$$X \subset \hat{Z} \subset Z \subset Y, \quad (\text{П.8.1})$$

$$\hat{Z} \subset X \subset Z \subset Y. \quad (\text{П.8.2})$$

Введем в рассмотрение пространство Y^* , сопряженное с Y относительно скалярного умножения в Z ; примем, что $Y^{**} = Y$, где Y^{**} - пространство, сопряженное с Y^* относительно того же скалярного умножения в Z .

2°. Выберем последовательность конечномерных подпространств $\{X_n\}$, полную в $X \cap Y^*$, и пусть

$$(y_{n1}, y_{n2}, \dots, y_{nr}), \quad n=1, 2, \dots \quad (\text{П.8.3})$$

- базис в X_n . Будем считать, что при любом n элементы (П.8.3) ортонормированы в \hat{Z} .

Обозначим через \hat{A} сужение оператора A на $Y^* \cap X$ и допустим, что \hat{A} - положительно определенный оператор из Y^* в Y (см. Бирман [1], а также книгу автора [8]). Допустим еще, что соответствующее энергетическое пространство совпадает с \hat{Z} - энергетическим пространством оператора \hat{A} .

Положим $Y_n = AX_n$. Проекторы Q_n определим соотношениями

$$\forall f \in Y; Q_n f \in Y_n, (f - Q_n f, y_{nj}) = 0, \quad j=1, 2, \dots, r; \quad (\text{П.8.4})$$

левая часть уравнения (П.8.4) рассматривается как значение функционала $(f - Q_n f) \in Y = Y^{**}$ на элементе $y_{nj} \in Y^*$. Как и в § 7,

доказывается, что Q_n есть проектор из Y на Y_n .

Оценим норму $\|Q_n\|_Y$, предполагая, что выполнено соотношение (П.8.1). Приближенное решение имеет вид

$$x_n^{(n)} = \sum_{k=1}^r a_k^{(n)} y_{nk};$$

при этом $A x_n^{(n)} = Q_n f$ и, следовательно,

$$\|Q_n f\|_Y^2 = \sum_{j, k=1}^N (A y_{nk}, A y_{nj})_Y a_k^{(n)} \bar{a}_{jn}^{(n)} \leq \Lambda_n^{(N)} \sum_{k=1}^N |a_k^{(n)}|^2; \quad (\text{П.8.5})$$

здесь, как и в § 7, $\Lambda_n^{(N)}$ есть наибольшее собственное число матрицы чисел $(A y_{nk}, A y_{nj})_Y$; $j, k = 1, 2, \dots, n$. Элементы y_{nk} ортонормированы в Σ , поэтому

$$\sum_{k=1}^N |a_k^{(n)}|^2 = \|x_*^{(n)}\|_{\Sigma}^2.$$

Нетрудно убедиться, что x_* есть решение задачи $\tilde{A}x = f$, а $x_*^{(n)}$ — приближенное решение этой задачи по Рунту. В книге автора [3] показано, что $|x_*^{(n)}| \leq |x_*|$; отсюда

$$\sum_{k=1}^N |a_k^{(n)}|^2 \leq |x_*| \leq C \|x_*\|_X = C \|A^{-1} f\|_X \leq C \|f\|_Y.$$

Теперь из (П.8.5) следует, что $\|Q_n\|_Y \leq C \sqrt{\Lambda_n^{(N)}}$ и

$$\|x_* - x_*^{(n)}\|_X \leq C \sqrt{\Lambda_n^{(N)}} \varepsilon_n(x_*). \quad (\text{П.8.6})$$

3°. Рассмотрим теперь случай (П.8.2). Неравенство (П.8.5) остается справедливым. Для произвольных элементов $h \in Y$ и $g \in Y'$ имеет место неравенство $|(h, g)| \leq \|h\|_Y \cdot \|g\|_{Y'}$. Используя уравнение (П.8.4) и ограниченность оператора \tilde{A}^{-1} , обратного положительно определенному оператору \tilde{A} , можно утверждать, что

$$\begin{aligned} \sum_{k=1}^N |a_k^{(n)}|^2 &= |x_*^{(n)}|^2 \leq |x_*|^2 = (\tilde{A}x_*, x_*) = (f, x_*) \leq \\ &\leq \|f\|_Y \cdot \|x_*\|_{Y'} = \|f\|_Y \cdot \|\tilde{A}^{-1} f\|_{Y'} \leq C \|f\|_Y^2. \end{aligned}$$

Теперь из (П.8.5) вытекает, что $\|Q_n f\|_Y \leq C \sqrt{\Lambda_n^{(N)}} \|f\|_Y$ и, следовательно, $\|Q_n\|_Y \leq C \sqrt{\Lambda_n^{(N)}}$. Отсюда окончательно

$$\|x_* - x_*^{(n)}\|_X \leq C \sqrt{\Lambda_n^{(N)}} \varepsilon_n(x_*). \quad (\text{П.8.7})$$

4°. Обратимся к задаче (П.4.1-2); сохраним все сделанные и предшествующих параграфах допущения о данных и решении этой задачи. Возьмем целое число l , $0 \leq l \leq 2s$, и обозначим $\tilde{s} = 2s - l$. Введем пространства $X = W_2^{(\tilde{s})}(\Omega)$, $Y = W_2^l(\Omega)$, $Z = L_2(\Omega)$

Тогда $X \subset Z \subset Y$; как и выше, нслик сверху означает, что функции из X удовлетворяют условиям (П.4.2). В данном случае $Z \sim \dot{W}_2^{(3)}(\Omega)$, поэтому расположения пространств (П.8.1) и (П.8.2) соответствуют значениям $0 < \nu < \delta$ и $\delta < \nu < 2\delta$. Чтобы доказать, что наша задача соответствует схеме п.2⁰ настоящего параграфа, достаточно установить справедливость следующих утверждений: а) операторы \hat{A} и \hat{A}^{-1} ограничены, б) оператор \hat{A} положительно определен, в) энергетическое пространство оператора \hat{A} совпадает с Z . Ограниченность оператора \hat{A} очевидна. Далее, в рассматриваемом нами случае $Y^* = X$; отсюда вытекает, что $\hat{A} = A$; нетрудно также убедиться, что A - симметричный и положительный оператор из $X = Y^*$ в Y , и что его энергетическое пространство совпадает с Z . Далее, если будет доказано, что оператор \hat{A}^{-1} ограничен, то положительный оператор \hat{A} будет положительно определенным. Окончательно, нам будет достаточно доказать ограниченность оператора \hat{A}^{-1} .

По определению, A есть расширение оператора \hat{A} , определяемое соотношением

$$(Ax, y)_Z = (x, \hat{A}^* y)_Z = (x, \hat{A} y)_Z; \quad (\text{П.8.8})$$

$$x \in \dot{W}_2^{(3)}(\Omega), \quad y \in W_2^{(2\delta)}(\Omega).$$

Оператор \hat{A} положительно определен в $Z = L_2(\Omega)$, поэтому он имеет только тривиальный ноль. Докажем, что у расширенного оператора A нет других нулей.

Пусть $Ax_0 = 0$. Тогда $(x_0, \hat{A}y)_Z = 0$. Это значит, что x_0 , рассматриваемый как функционал в $Z = L_2(\Omega)$, аннулируется на множестве значений оператора \hat{A} , т.е. на всем пространстве $L_2(\Omega)$. Но тогда $x_0 = 0$, оператор A не имеет нетривиальных нулей и существует обратный оператор A^{-1} . Далее, оператор \hat{A} очевидно нетеров, и его индекс равен нулю. (Как доказано в статье автора [9] (см. также [15]), оператор A также нетеров, и его индекс равен индексу оператора \hat{A} , т.е. равен нулю. Только что было показано, что подпространство нулей оператора A имеет нулевую размерность, а тогда нулевую размерность и подпространство нулей сопряженного оператора A^* . Отсюда следует, что оператор A^{-1} определен на всем пространстве $W_2^{(3-2\delta)}$. В то же время он замкнут как обратный нетерову (и, следовательно, замкнутому)

оператору A . Из сказанного следует, что оператор A^{-1} ограничен. 5° . Как и в § 7, имеем $x_n^{(n)}(t) = y^0(t)P_n(t)$, где $P_n(t)$ - полином степени $\leq n$ в R_m , и

$$\Lambda_n^{(n)} = \max_{\|t\|=1} \left\| \sum_{k=1}^n \tau_k A y_{nk} \right\|_{(-l)}^2 = \max_{\|t\|=1} \|A \varphi_n(\cdot, \tau)\|_{(-l)}^2.$$

Оператор A ограничен, поэтому

$$\|A \varphi_n(\cdot, \tau)\|_{(-l)} \leq C \|\varphi_n(\cdot, \tau)\|_{(2\delta)} = C \|\varphi_n(\cdot, \tau)\|_{(\delta)}.$$

Рассмотрим сначала случай $\tilde{\delta} > \delta$, или $l < \delta$, что соответствует расположению (П.8.1). По неравенству (П.6.6)

$$\|\varphi_n(\cdot, \tau)\|_{(\tilde{\delta})} \leq C n^{2(\tilde{\delta}-\delta)} \|\varphi_n(\cdot, \tau)\|_{(\delta)} \leq C n^{2(\tilde{\delta}-\delta)}. \quad (\text{П.8.9})$$

Далее, по теореме П.2.2,

$$\varepsilon_n(x_*) = \inf_{x \in X_n} \|x_* - x^{(n)}\|_{(\tilde{\delta})} \leq C n^{-2\tilde{\delta}} \omega_{L_2}^{(n)}\left(x_*, \frac{1}{n}\right),$$

и из формулы (П.8.6) следует

$$\delta < \tilde{\delta} \leq 2\delta, \quad \rho_n = \|x_* - x_*^{(n)}\|_{(\tilde{\delta})} \leq C n^{-2\tilde{\delta}+2\delta} \omega_{L_2}^{(n)}\left(x_*, \frac{1}{n}\right). \quad (\text{П.8.10})$$

Теперь рассмотрим случай $0 \leq \tilde{\delta} < \delta$. Неравенство (П.8.9) заменяется более простым:

$$\|\varphi_n(\cdot, \tau)\|_{(\tilde{\delta})} \leq C \|\varphi_n(\cdot, \tau)\|_{(\delta)} = C,$$

отсюда следует оценка

$$0 \leq \tilde{\delta} < \delta, \quad \rho_n = \|x_* - x_*^{(n)}\|_{(\tilde{\delta})} \leq C n^{-2\tilde{\delta}} \omega_{L_2}^{(n)}\left(x_*, \frac{1}{n}\right). \quad (\text{П.8.11})$$

Оценка (П.8.10) точнее, а оценка (П.8.11) совпадает с соответствующей оценкой статьи Шапошниковой [1].

Заметим, что при $\rho = 2$ и $\tilde{\delta} > \delta$ оценки Альина [5] совпадают с соответствующими оценками §§ 7, 8; при $\rho = 2$ и $\tilde{\delta} < \delta$ оценки Альина несколько хуже, чем соответствующие оценки § 4 ($\tilde{\delta} = \delta$)

и § 8 ($\delta < \delta$). Шапошникова показала в [3], что указанные здесь оценки Ильина могут быть улучшены.

Некоторые вопросы оценки погрешности метода Рунге в банаховых пространствах исследованы в статье Шапошниковой [4].

§ 9. Метод Бубнова - Галеркина

1°. Основы метода Бубнова-Галеркина, а также более общего метода Галеркина-Петрова изложены в книгах автора [2] и [4]; будем считать, что эти основы читателю известны.

В данном параграфе мы будем заниматься оценкой погрешности аппроксимации метода Бубнова-Галеркина, впервые, по-видимому, установленная Красносельским [1] в весьма общей ситуации - для нелинейных уравнений. В разных вариантах эта оценка изложена с подробными доказательствами в книге Красносельского и др. [1]; там же приведены необходимые литературные указания. Здесь мы ограничимся тем, что приведем (без доказательства) названную оценку в простейшем случае. Пусть A_0 - самосопряженный положительно определенный оператор в некотором гильбертовом пространстве H , K - линейный в H оператор, такой, что $\mathcal{D}(K) = \mathcal{D}(A_0)$ и $T = A_0^{-1}K$ - оператор, вполне непрерывный в энергетическом пространстве. Пусть $A = A_0 + K$, и уравнение $Ax = 0$ имеет только тривиальное решение. Тогда неоднородное уравнение $Ax = f$, $\forall f \in H$, имеет одно и только одно обобщенное решение (им является решение уравнения $x + Tx = A_0^{-1}f$). Обозначим это решение через x_* . В этих условиях (см. книгу автора [2]) при достаточно больших n существует одно и только одно приближенное решение $x_*^{(n)}$ по Бубнову - Галеркину, и $\|x_* - x_*^{(n)}\|_{n \rightarrow \infty} \rightarrow 0$; $\|\cdot\|$ - норма в H_0 . Упомянутая выше оценка погрешности имеет вид

$$\varepsilon_n(x_*) \leq \|x_* - x_*^{(n)}\| \leq C \delta_n(x_*); \quad (\text{п.9.1})$$

здесь $\varepsilon_n(x_*)$ имеет тот же смысл, что и в предшествующих параграфах настоящей главы.

Работы Вайникко [1, 3, 4, 9] содержит оценки аппроксимации для метода Бубнова-Галеркина в случае, когда координатные элементы суть собственные элементы самосопряженного оператора, обратного (т.е., имеющего общую область определения) с A_0 . Работы того же автора [5, 6, 8] посвящены оценкам погрешности аппроксимации метода Бубнова - Галеркина в проблеме собственных значений

Оценкам погрешности аппроксимации для метода Бубнова-Галеркина посвящены работы Джишкарлани [1, 4-7, 9].

2°. Приведем некоторые результаты перечисленных выше работ. В статье Джишкарлани [4] уравнение

$$A_0 x + Kx = f \quad (1.9.2)$$

рассматривают в предположениях, упомянутых в начале данного параграфа и при дополнительном предположении, что в пространстве H_0 вполне непрерывен оператор KA_0^{-1} . Пусть B - оператор, сходный с A_0 ; это значит, что B - самосопряженный положительно определенный оператор и $\mathcal{D}(B) = \mathcal{D}(A_0)$. Допустим еще, что B (a , следовательно, и A_0) имеет дискретный спектр; пусть ω_k и v_k - собственные числа и собственные элементы оператора B . Примем элементы v_k за координатные и обозначим через δ_n невязку приближенного решения $x_*^{(n)}: \delta_n = A_0 x_*^{(n)} + Kx_*^{(n)} - f$; допустим, что $\delta_n \in \mathcal{D}(B^v)$, $v \geq 0$; при $v = 0$ это допущение очевидно выполняется. Доказывается, что $\|\delta_n\|_H = O(1)$. В статье получены оценки

$$\|A_0(x_* - x_*^{(n)})\|_H = \frac{\|T_1\|_H \cdot \|B^v \delta_n\|}{\omega_{n+1}^v}; \quad (1.9.3)$$

$$|x_* - x_*^{(n)}| \leq \omega_{n+1}^{-v-1/2} [\|T_1\|_H \cdot \|T_1\|_H \cdot \|A_0^{-1} B\|_H]^{1/2} \|B^v \delta_n\|. \quad (1.9.4)$$

Здесь $|\cdot|$ - норма в H_0 - энергетическом пространстве оператора A_0 , I - тождественный оператор; $T = (I + A_0^{-1} K)^{-1}$, $T_1 = (I + KA_0^{-1})^{-1}$.

Вайникко в [9] рассмотрел более общее уравнение вида $Ax = f$, где A оператор, действующий из одного банахова пространства в другое. Мы изложим здесь результаты упомянутой статьи, относящиеся к уравнению (1.9.2). По-прежнему пусть B - оператор, сходный с A_0 , и пусть P_λ разложение единицы для оператора B . Обозначим $H_\lambda = P_\lambda H$. Приближенным решением уравнения (1.9.2) назовем на этот раз решение уравнения

$$P_\lambda(Ax_\lambda - f) = 0, \quad A = A_0 + K. \quad (1.9.5)$$

Допустим, что при некотором α_0 , $0 \leq \alpha_0 \leq 1$, оператор

$\Lambda^{d_0} \mathcal{K} \Lambda^{d_0}$ вполне непрерывен в H . Пусть T — замыкание этого оператора в H ; допустим, что уравнение $x + Tx = 0$ имеет только тривиальное решение. Пусть еще точное решение $x_* \in \mathcal{D}(B^{d_1})$. Тогда существует такое λ_0 , что при $\lambda \geq \lambda_0$ уравнение (П.9.5) имеет одно и только одно решение, и при $d_0 \leq d_1 \leq d_2$ верна оценка

$$\varepsilon_\lambda^{(d)} \leq \|B^{d_1}(x_* - x_\lambda)\| \leq C \varepsilon_\lambda^{(d)};$$

$$\varepsilon_\lambda^{(d)} = \|B^{d_1} x_* - P_\lambda B^{d_1} x_*\| = O(\lambda^{d_1 - d_2}) \quad (П.9.6)$$

3°. Рассмотрим краевую задачу

$$\sum_{|\alpha|, |\beta|=0}^{\delta} (-1)^{|\alpha|} \mathcal{D}^{d_1} (A_{\alpha\beta} \mathcal{D}^\beta x) = f(t), \quad t \in \Omega; \quad (П.9.7)$$

$$\mathcal{D}^\gamma x|_{\partial\Omega} = 0, \quad 0 \leq |\gamma| \leq \delta - 1.$$

допустим, что оператор этой задачи, который мы обозначим через A , можно представить в виде $A = A_0 + \mathcal{K}$, где оператор A_0 и положительно определен в $L_2(\Omega)$ и оператор $T = A_0^{-1} \mathcal{K}$ вполне непрерывен в H_0 — энергетическом пространстве оператора A_0 . Допустим далее, что граница области Ω достаточно гладкая, а задача (П.9.7) имеет не более одного решения; тогда решение этой задачи существует (может быть, как обобщенное). Примем, что это решение $x \in W_2^{(\nu)}(\Omega)$, где ν достаточно велико. В соответствии с частным случаем вложенной схемы Канторовича, описанным в п.2^с § 7, положим

$$X = W_2^{(\delta)}(\Omega), \quad Y = W_2^{(\nu)}(\Omega), \quad Z = L_2(\Omega).$$

Далее, пусть $\ell \geq 0$, $\tilde{\delta} = 2\delta + \ell$ (см. п.3^о § 7). Введем те же координатные функции, что и в п.2^с § 7, и определим проекторы тем же соотношением (П.7.6). Тогда остаются в силе все рассуждения § 7, и мы находим, что оценка (П.7.10) справедлива и для приближенного решения задачи (П.9.7), полученного по методу Гробо-ва — Галеркина при указанных выше координатных функциях.

Точно также нетрудно убедиться, что для задачи (П.9.7) верны оценки младших производных, полученные в § 8.

§ 10. Разностные методы.

Потребность аппроксимации разностных методов довольно хоро-

шо освещена в монографической литературе, и мы ограничимся здесь самыми краткими указаниями. Более подробно изложение можно найти, например, в монографиях Вазова и Форсайта [I], Бабушки, Витасака, Прагера [I], Марчука [I], Годунова и Рябенского [I].

Г°. Рассмотрим задачу Дирихле для уравнения Лапласа в двумерной области Ω с достаточно гладкой границей:

$$\Delta x = \frac{\partial^2 x}{\partial t_1^2} + \frac{\partial^2 x}{\partial t_2^2} = 0, \quad t = (t_1, t_2) \in \Omega, \quad x|_{\partial\Omega} = \varphi. \quad (\text{П.Ю.Г})$$

Построим квадратную сетку с шагом h . Обозначим через Ω_1 замкнутое подмножество Ω , обладающее тем свойством, что если $t \in \Omega_1$, то при любых числовых значениях $\xi, |\xi| \leq 1$, будет $(t_1 + \xi h, t_2) \in \bar{\Omega}$ и $(t_1, t_2 + \xi h) \in \bar{\Omega}$. В точках множества Ω_1 уравнение Лапласа аппроксимируется разностным уравнением

$$X(t_1, t_2) = \frac{1}{4} [X(t_1 + h, t_2) + X(t_1 - h, t_2) + X(t_1, t_2 + h) + X(t_1, t_2 - h)]; \quad (\text{П.Ю.2})$$

в точках множества $\Omega \setminus \Omega_1$ для функции X строится некоторая линейная интерполяционная формула, позволяющая имитировать краевое условие.

Узлы сетки, лежащие в Ω_1 , называются внутренними, а узлы, лежащие в $\Omega \setminus \Omega_1$ - граничными. Во внутренних узлах запишем уравнение (П.Ю.2), в граничных - упомянутую выше интерполяционную формулу. В результате получим некоторую линейно-алгебраическую систему (сеточную систему), о которой известно, что она имеет единственное решение X_h .

Доказывается (см. Вазов и Форсайт [I]), что если точное решение задачи (П.Ю.Г) $X_* \in C^{(k)}(\bar{\Omega})$ и X_h - функция, определенная на узлах сетки, лежащих в $\bar{\Omega}$, и совпадающая там с функцией X_* , то

$$\|X_h - X_*\|_{L_\infty} \leq \left(\frac{1}{12} M_k \sup_{\bar{\Omega}} Q + M_2 \right) h^2. \quad (\text{П.Ю.3})$$

Здесь Q - некоторая функция, непрерывная и неотрицательная в $\bar{\Omega}$, M_k - верхняя граница модулей частных производных от X_* порядка k .

Близкие результаты содержатся в той же монографии для общего эллиптического уравнения 2-го порядка в области Ω любой размерности и для более общих краевых условий.

В монографии Бабушки, Битаева и Прагера [1] даны оценки погрешности аппроксимации разностного метода для бигармонического уравнения на плоскости.

2°. Рябенкий и Филиппов доказали весьма интересную общую теорему об оценке погрешности аппроксимации разностных методов. Эта теорема подробно доказывается в монографиях Годунова и Рябенкого [1] и Марчука [1]. Рассмотрим вычислительную задачу вида

$$Lx = f, \quad lx = g. \quad (\text{П.10.4})$$

Здесь x и f - функции, искомая и заданная, определенные в области $\Omega \subset \mathbb{R}^m$, g - функция, определенная на $\partial\Omega$, L и l - линейные операторы. Пусть задача (П.10.4) заменена разностной задачей

$$L_h x^h = f^h, \quad l_h x^h = g^h, \quad (\text{П.10.5})$$

в которой x^h, f^h, g^h - сеточные функции, преобразующие, в общем случае, различным пространствам сеточных функций, так что $x^h \in X_h, f^h \in \mathcal{F}_h, g^h \in \mathcal{G}_h$.

Пусть $(\cdot)_h$ - оператор, преобразующий функцию, заданную в Ω или на $\partial\Omega$, в сеточную функцию, заданную на соответствующем множестве узлов. Теорема Рябенкого - Филиппова состоит в следующем:

Пусть справедливы аппроксимационные неравенства

$$\begin{aligned} \|(Lx)_h - L_h x^h\|_{X_h} &\leq Ch^k, \quad \|(lx)_h - l_h x^h\|_{\mathcal{G}_h} \leq Ch^k, \\ \|(f)_h - f^h\|_{\mathcal{F}_h} &\leq Ch^k, \quad \|(g)_h - g^h\|_{\mathcal{G}_h} \leq Ch^k \end{aligned} \quad (\text{П.10.6})$$

и пусть обратный оператор задачи (П.10.5) ограничен независимо от h , так что

$$\|x^h\|_{X_h} \leq C \{ \|f^h\|_{\mathcal{F}_h} + \|g^h\|_{\mathcal{G}_h} \}. \quad (\text{П.10.7})$$

Тогда справедлива оценка погрешности аппроксимации

$$\|(x)_h - x^h\| \leq Ch^k. \quad (\text{П.10.8})$$

§ II. Метод коллокации.

1°. Этот метод был впервые предложен Кантсоровичем и его

статье [1]. Сущность метода такова. Пусть требуется приближенно решить линейное уравнение

$$Ax = f; \quad x \in X, \quad f \in Y, \quad (\text{П. II.1})$$

(требование линейности несущественно), в котором X и Y - банаховы пространства. A - оператор, действующий из X в Y причём элементы пространства Y образуют подмножество множества функций, непрерывных на некотором компакте $K \subset R_m$. Выберем некоторую последовательность конечномерных подпространств $\{X_n\}$, полных в X , и положим $\dim X_n = N(n) = N$, $Y_n = AX_n$.

Пусть существует оператор A^{-1} , определенный на всем пространстве Y , тогда $\dim Y_n = N$ и последовательность $\{Y_n\}$ полна в Y . Далее, если $\{y_{nk}\}$, $k = 1, 2, \dots, N$ - базис в X_n , и $\psi_{nk}(t) = Ay_{nk}, t \in K$, то $\{\psi_{nk}\}$ есть базис в Y_n .

Выберем в K точки $t_j = t_j^{(n)}, j = 1, 2, \dots, N$, называемые узлами коллокации. Приближенное решение ищем как элемент подпространства X_n ,

$$x^{(n)} = \sum_{k=1}^N a_k^{(n)} y_{nk} \quad (\text{П. II.2})$$

и определяем коэффициенты $a_k^{(n)}$ из условия, чтобы уравнение (П. II.1) выполнялось в узлах коллокации:

$$\sum_{k=1}^N a_k^{(n)} \psi_{nk}(t_j) = f(t_j). \quad (\text{П. II.3})$$

2°. В ряде работ исследуется погрешность аппроксимации рассматриваемого здесь метода. Приведем здесь важнейшие результаты некоторых из этих работ. Сразу же отметим, что некоторые оценки упомянутой погрешности, а также дополнительная библиография по этому вопросу, приведены в книгах Канторовича и Акилова [1], Красильского и др. [1].

В работах Карпиловой [1 - 3] изучена погрешность аппроксимации коллокационного метода для обыкновенного дифференциального уравнения порядка 2 при условиях первой краевой задачи. Пусть уравнение задано на промежутке $a \leq t \leq b$, и пусть точное решение $x_* \in C^{(r,2)}[a,b], 0 < r \leq 1$. Приближенное решение $x_*^{(n)}$ строится в виде полинома степени $k+2s$, удовлетворяющего краевой условиям задачи. Пусть коэффициенты и свободный член дифференциального уравнения принадлежат классу $C^{(r,2)}[a,b]$, и пусть задача

имеет не более одного решения. Если в качестве узлов коллокации выбраны узлы Чебышева или Гаусса, то соответственно

$$\|x_* - x_*^{(n)}\|_{C(\sigma)} \leq C \frac{\ln n}{n^{\alpha+d}} \quad (\text{П. II.4})$$

$$\|x_* - x_*^{(n)}\|_{C(\omega)} \leq \frac{C}{n^{\alpha-d/2}} \quad (\text{П. II.5})$$

в цитированных работах Карпиловской рассмотрена также первая краевая задача для уравнения

$$\Delta x + \lambda p(t)x = f(t) \quad (\text{П. II.6})$$

в квадрате $R: 0 \leq t_1, t_2 \leq \pi$. За координатные функции взяты собственные функции той же задачи для уравнения Лапласа:

$$y_{ik}(t) = \sin it_1 \sin kt_2; \quad 0 < i \leq m, \quad 0 < k \leq n;$$

узлы коллокации - точки $(\frac{2i-1}{2m+1}\pi, \frac{2k-1}{2n+1}\pi)$. Если λ отлично от собственных чисел рассматриваемой задачи, и $p, f \in \text{Lip}_\alpha(R)$, то $(x_*^{(m,n)})$ - приближенное решение

$$\|\Delta(x_* - x_*^{(m,n)})\|_{C(R)} \leq C \ln^2 m \ln^2 n \left(\frac{1}{m^2} + \frac{1}{n^2} \right). \quad (\text{П. II.7})$$

в статьях Карпиловской [2, 3] рассмотрены и некоторые другие задачи.

3°. В работе Габдулхасова [1] рассматривается интегральное уравнение Фредгольма

$$Ax(t) = \mu x(t) - \frac{1}{2\pi} \int_0^{2\pi} K(t,s)x(s)ds = f(t) \quad (\text{П. II.8})$$

в пространстве 2π -периодических функций с C -нормой. Приближенное решение строится в виде кусочно-линейной функции с угловыми точками $t_j = t_j^{(n)} = 2\pi j/n$; эти же точки суть узлы коллокации. Доказывается следующее утверждение. Пусть μ отлично от собственных чисел уравнения (П. II.3), а K таково, что

$$\nu := \omega_\nu(K, \frac{2\pi}{n}) \|A^{-1}\| < 1. \quad (\text{П. II.9})$$

Тогда алгебраическая система метода коллокации имеет единственное решение. При этом, если $x_*(t)$ - точное, а $x_*^{(n)}(t)$ - приближенное решение уравнения (П. II.8), то

$$\|x_n - x_n^{(n)}\|_C \leq (1 + P_A / (1 - \tau)) \omega(x_n, \frac{2\pi}{n}). \quad (\text{П. II. IO})$$

Здесь $P_A = \|A\| \cdot \|A^{-1}\|$ - число обусловленности оператора A , а ω в двух последних формулах означает модуль непрерывности.

В той же статье Габдулхаева [I] содержатся и некоторые результаты, относящиеся к сингулярным одномерным интегральным уравнениям.

4°. В книге Вильбермана и Пресдорфа [I] (см. также книгу Михлина и Пресдорфа [I]) метод коллокации применен к сингулярному одномерному интегральному уравнению вида

$$a(t)x(t) + \frac{b(t)}{\pi i} \int_{\Gamma} \frac{x(s)}{s-t} ds + Tx = f(t). \quad (\text{П. II. II})$$

Здесь Γ - окружность $|t| = 1$, $a(t)$ и $b(t)$ непрерывны и символ не вырождается, так что $a^2(t) - b^2(t) \neq 0$, $t \in \Gamma$; T - оператор вполне непрерывный в $L_p(\Gamma)$, $1 < p < \infty$. Координатные функции-тригонометрические полиномы порядка n относительно $\theta = \arg t$, узлы коллокации, обладают с точками $t_j^{(n)} = \exp(\frac{2\pi i j}{n})$. Пусть

$$\tilde{a}(t) = a(t) + b(t), \quad \tilde{b}(t) = a(t) - b(t); \quad \text{ind } \tilde{a}(t) = \text{ind } \tilde{b}(t); \\ a(t), b(t) \in C^{(\mu, \nu)}(\Gamma); \quad T: L_p(\Gamma) \longrightarrow C^{(\mu, \nu)}(\Gamma).$$

Тогда

$$\|x_n - x_n^{(n)}\|_p = O(n^{-\lambda}), \quad \lambda = \min(\mu, \nu). \quad (\text{П. II. I2})$$

Некоторые результаты получены и для уравнений вида (П. II. II) с вырождающимся символом.

Пресдорф и Шмидт [I] рассмотрели метод коллокации для уравнения (П. II. II) при $T = 0$ с узлами $t_j^{(n)} = \exp(\frac{2\pi i j}{n})$ и с кусочно-линейными координатными функциями, угловые точки которых совпадают с узлами коллокации. Основным результатом упомянутой статьи является следующий: пусть коэффициенты $a(t)$ и $b(t)$ непрерывны. Для того, чтобы данный метод коллокации сходился в $L_p(\Gamma)$ необходимо и достаточно, чтобы

$$a(t) + \lambda b(t) \neq 0; \quad \forall t \in \Gamma, \quad \forall \lambda \in [-1, +1]. \quad (\text{П. II. I3})$$

Если при этом $a, b \in \text{Lip}_\mu(\Gamma)$, $f \in \text{Lip}_\lambda(\Gamma)$, то

$$\|x_* - x_*^{(n)}\|_{C^{\nu}(\Gamma)} = O(n^{-\nu} \ln n), \quad \nu = \min(\lambda, \mu). \quad (\text{п. II.14})$$

Некоторые результаты получены в цитированной статье и для разрывных символов.

5°. Ряд результатов по оценке погрешности аппроксимации для уравнений Фредгольма и одномерных сингулярных интегральных уравнений содержится в книге Габдулхаева [3].

ГЛАВА III

Погрешность искажения

§ I. Погрешность искажения линейного свободного вычислительного процесса

1^o. Погрешность искажения вычислительных процессов исследована в ряде работ автора и его учеников; первая из работ автора этого направления [4] опубликована в 1960 г.; работы периода 1960 - 65 гг. суммированы в книге автора [8]; эти работы тесно связаны понятием, тогда же введенным автором, об устойчивости вычислительного процесса по отношению к искажению. В книге [8] содержатся и не публиковавшиеся ранее результаты, в основном касающиеся нелинейных вычислительных процессов; эти процессы будут рассмотрены ниже, в главах УД-IX. Ряд результатов был позднее получен для погрешности искажения метода конечных элементов (см. работы автора [13, 14]); результаты этого плана будут рассмотрены в гл. V.

В данной главе мы изложим результаты, относящиеся к искажению линейных вычислительных процессов и полученные частично в упомянутых выше более старых работах, частично в более новой работе автора [33].

2^o. Рассмотрим линейный свободный вычислительный процесс, состоящий в решении последовательности независимых уравнений

$$A_n x^{(n)} = f^{(n)}; \quad x^{(n)} \in X_n, \quad f^{(n)} \in Y_n, \quad n=0, 1, 2, \dots \quad (\text{Ш. I. 1})$$

Здесь X_n и Y_n - банаховы пространства, A_n - замкнутый линейный оператор, отображающий X_n на Y_n . Мы принимаем, что операторы A_n ограниченно обратимы, так что каждый из операторов A_n существует, ограничен и определен на всем пространстве Y_n .

Искаженный вычислительный процесс, который естественно так же считать линейным, имеет вид

$$(A_n + \Gamma_n) x^{(n)} = f^{(n)} + \delta^{(n)}; \quad n=0, 1, 2, \dots \quad (\text{Ш. I. 2})$$

Можно считать, что оператор Γ_n - искажение оператора A_n , в том или ином смысле мал по сравнению с A_n . Примеч, что величина $\|A_n^{-1} \Gamma_n\|$ может быть сделана сколь угодно малой. Норму элемента $\delta^{(n)}$ также можно считать сколь угодно малой, однако в пределах линейной теории это допущение не играет существенной роли.

для последующих рассуждений.

Оценим погрешность искажения процесса (Ш. I. I). Зададим число β , $0 < \beta < 1$, и потребуем, чтобы $\|A_n^{-1} \Gamma_n\| \leq \beta$. При таком условии оператор $A_n + \Gamma_n$ ограниченно обратим. Действительно (I_n - тождественный оператор в X_n) $A_n + \Gamma_n = A_n(I_n + A_n^{-1} \Gamma_n)$. Первый

множитель ограниченно обратим по предположению, второй - по известной теореме Банаха, и

$$(A_n + \Gamma_n)^{-1} = (I_n + A_n^{-1} \Gamma_n)^{-1} A_n^{-1}. \quad (\text{Ш. I. 3})$$

Решение уравнения (Ш. I. 2) существует и единственно:

$$x_*^{(n)} = (I_n + A_n^{-1} \Gamma_n)^{-1} A_n^{-1} (f^{(n)} + \delta^{(n)}).$$

Решение уравнения (Ш. I. I) $x_*^{(n)} = A_n^{-1} f^{(n)}$, поэтому

$$x_*^{(n)} - x_*^{(n)} = [(I_n + A_n^{-1} \Gamma_n)^{-1} - I_n] x_*^{(n)} + (I_n + A_n^{-1} \Gamma_n)^{-1} A_n^{-1} \delta^{(n)}.$$

Если I тождественный, а $I + \gamma$ ограниченно обратимый оператор в некотором банаховом пространстве, то $(I + \gamma)[(I + \gamma)^{-1} - I] = -\gamma$ и, следовательно $(I + \gamma)^{-1} - I = -(I + \gamma)^{-1} \gamma$. В частности,

$$(I_n + A_n^{-1} \Gamma_n)^{-1} - I_n = -(I_n + A_n^{-1} \Gamma_n)^{-1} A_n^{-1} \Gamma_n; \quad (\text{Ш. I. 5})$$

замечая еще, что $\|(I_n + A_n^{-1} \Gamma_n)^{-1}\| \leq (1 - \beta)^{-1}$, мы приходим к

исковой оценке погрешности искажения:

$$\hat{\xi}_n = \|x_*^{(n)} - x_*^{(n)}\| \leq \frac{1}{1 - \beta} [\|A_n^{-1} \Gamma_n x_*^{(n)}\| + \|A_n^{-1} \delta^{(n)}\|]. \quad (\text{Ш. I. 6})$$

Отсюда вытекает несколько более простая и более удобная для применения оценка:

$$\hat{\xi}_n \leq \frac{1}{1 - \beta} [\|A_n^{-1} \Gamma_n\| \cdot \|x_*^{(n)}\| + \|A_n^{-1} \delta^{(n)}\|]. \quad (\text{Ш. I. 7})$$

Формула (Ш. I. 7) позволяет оценить погрешность искажения процесса (Ш. I. I), если известны оценки точных решений $x_*^{(n)}$ уравнения (Ш. I. I).

§ 2. Устойчивость свободного процесса относительно искажения

Важнейшим и интересным условием устойчивости является ограниченность операторов A_n и Γ_n . Если A_n и Γ_n ограничены, то оператор $A_n + \Gamma_n$ ограниченно обратим, если $\|A_n^{-1} \Gamma_n\| < 1$. Если A_n и Γ_n ограничены, то оператор $A_n + \Gamma_n$ ограниченно обратим, если $\|A_n^{-1} \Gamma_n\| < 1$.

в том или ином смысле малы искажения Γ_n и $\delta^{(n)}$. Такие условия связаны с понятием устойчивости вычислительного процесса относительно искажения; это понятие сформулировано в книге автора [8] и несколько уточнено в его работах [17, 22].

Целесообразными оказываются два несколько различных определения устойчивости. Согласно первому из них мы называем процесс (Ш.1.1) устойчивым относительно погрешностей искажения в последовательности пар пространств (X_n, Y_n) , если существуют такие положительные числа ρ, q, ν , что из неравенства $\|A_n^{-1} \Gamma_n\| \leq \nu$ следует однозначная разрешимость уравнения (Ш.1.2) при любом n и неравенство

$$\|x_*^{(n)} - \alpha_*^{(n)}\| \leq \rho \|A_n^{-1} \Gamma_n\| + q \|\delta^{(n)}\|. \quad (\text{Ш.2.1})$$

Второе определение имеет смысл, если искажения Γ_n суть ограниченные операторы. В этом случае можно назвать процесс (Ш.1.1) устойчивым относительно искажения в последовательности пар пространств (X_n, Y_n) , если существуют такие положительные постоянные ρ, q, ν , что из неравенства $\|\Gamma_n\| \leq \nu$ вытекает однозначная разрешимость уравнения (Ш.1.2) при любом n и неравенство

$$\|x_*^{(n)} - \alpha_*^{(n)}\| \leq \rho \|\Gamma_n\| + q \|\delta^{(n)}\|. \quad (\text{Ш.2.2})$$

Ниже мы для краткости будем говорить, что вычислительный процесс устойчив (или неустойчив), если он устойчив (или неустойчив) относительно погрешностей искажения; там, где это не может вызвать недоразумений, мы будем опускать слова "в последовательности пар пространств (X_n, Y_n) ".

2°. Теорема Ш.2.1. Для того, чтобы вычислительный процесс (Ш.1.1) был устойчив в смысле первого определения, необходимо достаточно, чтобы

$$\|A_n^{-1}\| \leq C_1, \quad \|x_*^{(n)}\| \leq C_2; \quad C_1, C_2 = \text{const}. \quad (\text{Ш.2.3})$$

Достаточность условий (Ш.2.3) сразу вытекает из оценки (Ш.1.7): если указанные условия выполнены, то можно положить $\nu = \beta$, $0 < \beta < 1$, и из (Ш.1.6) следует неравенство (Ш.2.1) с постоянными

$$\rho = \frac{C_1}{1-\beta}, \quad q = \frac{C_1}{1-\beta}.$$

Необходимость. Пусть процесс (Ш.1.1) устойчив в смысле пе

вого определения. Допустим, что $\Gamma_n = 0$, так что покажутся только свободные члены уравнений (Ш.1.1). Неравенство (Ш.2.1) принимает вид $\|y^{(n)}\| \leq q \|\delta^{(n)}\|$, где $y^{(n)} = x_*^{(n)} - x_*^{(n)}$. Далее, в данном случае $A_n x_*^{(n)} = f^{(n)} + \delta^{(n)}$, поэтому $A_n y^{(n)} = \delta^{(n)}$ и, следовательно, $\|A_n y^{(n)}\| \geq q^{-1} \|y^{(n)}\|$. Это означает, что нормы $\|A_n^{-1}\|$ ограничены независимо от n ; можно положить $C_1 = q^{-1}$.

Чтобы доказать необходимость второго условия (Ш.2.3), допустим, что $\delta^{(n)} = 0$. Тогда, если только $\|A_n^{-1} \Gamma_n\| \leq \nu$, то

$$\|x_*^{(n)} - x_*^{(n)}\| \leq \rho \|A_n^{-1} \Gamma_n\|. \quad (\text{Ш.2.4})$$

Положим $\Gamma_n = \nu A_n$. Уравнение (Ш.1.2) принимает вид

$$(1+\nu) A_n x^{(n)} = f^{(n)}. \text{ Отсюда } x_*^{(n)} = (1+\nu)^{-1} x_*^{(n)} \text{ и из (Ш.2.4) следует, что } \|x_*^{(n)}\| \leq \rho(1+\nu) = \text{const}.$$

3°. Теорема Ш.2.2. (см. работы автора [7, 8]). Для того, чтобы процесс (Ш.1.1) был устойчив в смысле второго определения в последовательности (X_n, Y_n) , необходимо и достаточно, чтобы были выполнены следующие условия: 1) $\|A_n\| \leq C_1 = \text{const}$; 2) существует такая постоянная C_2 , что каков бы ни был линейный оператор B_n с единичной нормой, действующий из X_n в Y_n , справедливо неравенство

$$\|A_n^{-1} B_n x_*^{(n)}\| \leq C_2. \quad (\text{Ш.2.5})$$

Необходимость условия $\|A_n^{-1}\| \leq C_1$ доказывается так же, как и в теореме Ш.2.1. Докажем необходимость условия (Ш.2.5). Пусть $\delta^{(n)} = 0$. Положим $\Gamma_n = \varepsilon B_n$, $\varepsilon = \text{const} < \nu$. Тогда $(A_n + \Gamma_n) x_*^{(n)} = f^{(n)}$, отсюда $x_*^{(n)} + A_n^{-1} \Gamma_n x_*^{(n)} = x_*^{(n)}$ и

$$x_*^{(n)} - x_*^{(n)} + A_n^{-1} \Gamma_n (x_*^{(n)} - x_*^{(n)}) = -A_n^{-1} \Gamma_n x_*^{(n)}. \quad (\text{Ш.2.6})$$

Далее,

$$\|A_n^{-1} \Gamma_n x_*^{(n)}\| = \varepsilon \|A_n^{-1} B_n x_*^{(n)}\| = \|\Gamma_n\| \cdot \|A_n^{-1} B_n x_*^{(n)}\|;$$

$$\|A_n^{-1} \Gamma_n (x_*^{(n)} - x_*^{(n)})\| = \varepsilon \|A_n^{-1} B_n (x_*^{(n)} - x_*^{(n)})\| \leq$$

$$\leq \varepsilon \|A_n^{-1}\| \cdot \|x_*^{(n)} - x_*^{(n)}\| \leq \varepsilon C_1 \|x_*^{(n)} - x_*^{(n)}\|,$$

и из (III.2.6) следует

$$\|x_*^{(n)} - x_*^{(n)}\| \geq \frac{\|\Gamma_n\| \cdot \|A_n^{-1} B_n x_*^{(n)}\|}{1 + \varepsilon C_1}.$$

Сравнив это с неравенством устойчивости (III.2.2), которое в данном случае принимает вид $\|x_*^{(n)} - x_*^{(n)}\| \leq \rho \|\Gamma_n\|$, получаем

$$\|A_n^{-1} B_n x_*^{(n)}\| \leq \rho(1 + \varepsilon C_1).$$

Пологая здесь $\varepsilon \rightarrow 0$, находим, что условие (III.2.5) выполнено со значением $C_3 = \rho$.

Теперь докажем достаточность условий теоремы. Положим $\tau = \beta/C_1, \beta = \text{const}, 0 < \beta < 1$, и потребуем, чтобы $\|\Gamma_n\| \leq \tau$. Если $\Gamma_n = 0$, то $x_*^{(n)} = x_*^{(n)} + A_n^{-1} \delta^{(n)}$, $\|x_*^{(n)} - x_*^{(n)}\| \leq C_1 \|\delta^{(n)}\|$ и теорема доказана. Если же $\Gamma_n \neq 0$, то положим $B_n = \Gamma_n / \|\Gamma_n\|$ и $x_*^{(n)} = u^{(n)} + v^{(n)}$, где

$$(A_n + \Gamma_n) u^{(n)} = f^{(n)}, \quad (A_n + \Gamma_n) v^{(n)} = \delta^{(n)}. \quad (\text{III.2.7})$$

Имеем

$$\|v^{(n)}\| \leq \|(I_n + A_n^{-1} \Gamma_n)^{-1}\| \cdot \|A_n^{-1}\| \cdot \|\delta^{(n)}\| \leq \frac{C_1}{1 - \beta} \|\delta^{(n)}\|. \quad (\text{III.2.8})$$

По формуле (III.1.5) находим из первого уравнения (III.2.7)

$$\|u^{(n)} - x_*^{(n)}\| = \|(I_n + A_n^{-1} \Gamma_n)^{-1} A_n^{-1} \Gamma_n x_*^{(n)}\| \leq \frac{\|A_n^{-1} B_n x_*^{(n)}\| \|\Gamma_n\|}{1 - \beta} \leq \frac{C_2}{1 - \beta} \|\Gamma_n\|;$$

неравенство устойчивости (III.2.2) выполняется при значениях постоянных $p = (1 - \beta)^{-1} C_2, q = (1 - \beta)^{-1} C_1$.

4°. Следствие III.2.1. Если нормы $\|x_*^{(n)}\|$ ограничены в совокупности, то для устойчивости процесса (III.1.1) в смысле второго определения необходимо и достаточно, чтобы $\|A_n^{-1}\| \leq C_1 = \text{const}$.

Сведем частный случай, обычно встречающийся при реализации итерационных методов. Пусть X_n — подпространства пространства X , и пусть $x_*^{(n)} \xrightarrow{n \rightarrow \infty} x_*$; в этом случае назовем вычисли-

тельный процесс (Ш.1.1) сходиться. Для такого процесса величины $\|x_*^{(n)}\|$ ограничены, и для его устойчивости необходимо и достаточно, чтобы $\|A_n^{-1}\| \leq C_1 = \text{const}$.

Следствие Ш.2.2. (Суворов [1]). Для устойчивости процесса (Ш.1.1) в смысле второго определения необходимо и достаточно, чтобы

$$\|A_n^{-1}\| \leq C_1, \|A_n^{-1}\| \cdot \|x_*^{(n)}\| \leq C_3; C_1, C_3 = \text{const}. \quad (\text{Ш.2.9})$$

Достаточность условий (Ш.2.9) и необходимость первого из них очевидны; докажем необходимость второго условия.

Пусть процесс устойчив. В силу теоремы Ш.2.2, $\|A_n^{-1} B_n x_*^{(n)}\| \leq C_2$,

$\forall B_n, \|B_n\| = 1$. Допустим, что произведение $\|A_n^{-1}\| \cdot \|x_*^{(n)}\|$ неограничено. Перейдем к подпоследовательности, можно считать, что

$\|A_n^{-1}\| \cdot \|x_*^{(n)}\| \xrightarrow{n \rightarrow \infty} \infty$. Для каждого n найдем такой элемент $y^{(n)}$, чтобы $\|y^{(n)}\| = 1$ и $\|A_n^{-1} y^{(n)}\| \geq \frac{1}{2} \|A_n^{-1}\|$. Построим далее

линейный функционал ℓ_n , определенный на X_n , со свойствами $\|\ell_n\| = 1, \|\ell_n x_*^{(n)}\| = \|x_*^{(n)}\|$. Положим $B_n x_*^{(n)} = \ell_n x_*^{(n)} \cdot y^{(n)}$;

ясно, что $\|B_n\| = 1$. Имеем $A_n^{-1} B_n x_*^{(n)} = \|x_*^{(n)}\| \cdot A_n^{-1} y^{(n)}$ и

$\|A_n^{-1} B_n x_*^{(n)}\| \geq \frac{1}{2} \|A_n^{-1}\| \cdot \|x_*^{(n)}\| \xrightarrow{n \rightarrow \infty} \infty$, что противоречит предположению.

Из следствия Ш.2.2 сразу вытекает

Следствие Ш.2.3. Если нормы $\|A_n^{-1}\|$ положительного ограничены снизу, то для устойчивости процесса (Ш.1.1) в смысле второго определения необходимо и достаточно, чтобы $\|A_n^{-1}\| \leq C_1, \|x_*^{(n)}\| \leq C_2$;

$C_1, C_2 = \text{const}$. Если еще процесс (Ш.1.1) — сходится, то для его устойчивости необходимо и достаточно, чтобы $\|A_n^{-1}\| \leq C_1$.

Следствие Ш.2.3 доказано в книге автора [8] независимо от следствия Ш.2.2.

5°. Если $\sup_n \|A_n^{-1}\| = \infty$, то процесс (Ш.1.1) неустойчив в последовательности (X_n, Y_n) . Рассмотрим общий случай, когда

$\sup_n \|A_n^{-1}\|_{y_n \rightarrow x_n} \leq \infty$. Возьмем новое пространство \bar{X}_n и \bar{Y}_n , которые совпадают соответственно с X_n и Y_n как множества элементов, но члены в них описываются формулами

$$\| \cdot \|_{\bar{X}_n} = \frac{1}{\sqrt{\|A_n^{-1}\|}} \| \cdot \|_{X_n}; \quad \| \cdot \|_{\bar{Y}_n} = \sqrt{\|A_n^{-1}\|} \| \cdot \|_{Y_n}. \quad (\text{Ш.2.10})$$

Теорема Ш.2.3. Для того, чтобы процесс (Ш.1.1) был устойчивым в смысле второго определения в последовательности (\bar{X}_n, \bar{Y}_n) , необходимо и достаточно выполнение неравенства

$$\|x_*^{(n)}\|_{\bar{X}_n} \leq C = \text{const}. \quad (\text{Ш.2.11})$$

Доказательство сразу вытекает из соотношения (Ш.2.9), найденных для пространств X_n, Y_n , и из очевидного тождества

$$\|A_n^{-1}\|_{\bar{Y}_n \rightarrow \bar{X}_n} = 1.$$

Замечание Ш.2.1. Теорема Ш.2.3 остается в силе, если нормы в \bar{X}_n и \bar{Y}_n определить формулами

$$\| \cdot \|_{\bar{X}_n} = \gamma(n) \| \cdot \|_{X_n}; \quad \| \cdot \|_{\bar{Y}_n} = \frac{1}{\gamma(n)} \| \cdot \|_{Y_n}, \quad (\text{Ш.2.12})$$

где $\gamma(n)$ - любая положительная функция от n , такая, что

$$\|A_n^{-1}\|_{\bar{Y}_n \rightarrow \bar{X}_n}^{-1} \geq C\gamma^2(n), \quad C = \text{const}. \quad (\text{Ш.2.13})$$

§ 3. Устойчивость процесса Рунге

1°. Метод Рунге коротко описан выше, в § 1, гл. II. Обозначим соответственно через $a^{(n)}$ и $f^{(n)}$ векторы с составляющими $a_k^{(n)}$ и (f, y_{nk}) , $k = 1, 2, \dots, N$, а через M_n - матрицу элементов $[y_{nk}, y_{nj}]$; $k, j = 1, 2, \dots, N$ - матрицу Рунге, и запишем систему (П.1.4) в виде

$$M_n a^{(n)} = f^{(n)}; \quad n = 1, 2, \dots \quad (\text{Ш.3.1})$$

Будем рассматривать $a^{(n)}$ и $f^{(n)}$ как элементы пространства R_n , а M_n как оператор в R_n . Если применить к данному случаю общую схему § 1 гл. I, то $X = H_A, X_n = R_n, x^{(n)} = a^{(n)}$; соотношения

x_n задается системой (П.1.4) или, что то же, системой (Ш.3.1). Как и в операторе R_n , отображающий R_n в H_A , определяется по формуле (П.1.3).

2°. С методом Рунге связаны два вычислительных процесса: процесс (Ш.3.1) приводящий к вычислению вектора коэффициентов $a^{(n)}$, и процесс вычисления приближенных решений $x_*^{(n)}$. Если обозначить через a_* решение уравнения (Ш.1.1), то

$$x_*^{(n)} = P_n a_*^{(n)}; \quad (\text{Ш.3.2})$$

нмем в виду применить общие теоремы § 2 настоящей главы, запишем этот процесс иначе. Будем рассматривать R_n как оператор из R_n в $X_n = H_n^{(n)}$, тогда существует обратный оператор $S_n = R_n^{-1}$; $S_n x^{(n)} = a^{(n)}$. Подставив это в (Ш.1.2), получим процесс, определяющий приближенные решения:

$$M_n S_n x^{(n)} = f^{(n)}. \quad (\text{Ш.3.3})$$

Ниже в данном параграфе мы исследуем устойчивость обоих указанных здесь процессов в смысле второго определения § 2 настоящей главы.

5°. Пусть в некотором гильбертовом пространстве задана последовательность вида (П.1.5) Пусть, далее, G_n - матрица Грама элементов $y_{n1}, y_{n2}, \dots, y_{nN}$ и $\lambda_n^{(1)}$ - наименьшее собственное число этой матрицы. Последовательность (П.1.5) назовем *сильно минимальной* в H , если числа $\lambda_n^{(1)}$ ограничены снизу, положительным числом, не зависящим от n .

Понятие *сильной минимальности* Гел'фандина [1] для последовательностей вида (П.1.6); некоторые свойства таких последовательностей рассмотрены в книге автора [8].

Теорема Ш.3.1. Для того, чтобы процесс (Ш.3.1) вычисления коэффициентов Рунга был устойчив в последовательности (R_n, R_n) , необходимо и достаточно, чтобы координатная система была *сильно минимальной* в энергетическом пространстве рассматриваемой задачи.

Матрица Рунга совпадает с матрицей Грама (в энергетической метрике) координатных элементов $y_{n1}, y_{n2}, \dots, y_{nN}$. Эта матрица - симметричная и положительно определенная, и ее собственные числа положительны; пусть $\lambda_n^{(1)}$ - наименьшее из них. Обозначим через A_n оператор, порожденный матрицей M_n и действующий в R_n . В данном случае элемент $x_*^{(n)}$ (§ 2 гл.1) следует отождествить с $a_*^{(n)}$ - решением системы (Ш.3.1). Очевидно, $\|A_n^{-1}\| = 1/\lambda_n^{(1)}$. По теореме Ш.2.1, для устойчивости необходимо, чтобы $\|A_n^{-1}\| \leq C = \text{const}$, или $\lambda_n^{(1)} \geq 1/C$. Чтобы доказать достаточность условия теоремы, остается проверить (см. следствие Ш.2.1), что нормы $\|a_*^{(n)}\|$ ограничены. Процесс Рунга - сходящийся в H_A , поэтому $\|x_*^{(n)}\| \leq C = \text{const}$. Далее,

$$|x_*^{(n)}|^2 = \left[\sum_{k=1}^N a_k^{(n)} y_{nk}, \sum_{j=1}^N a_j^{(n)} y_{nj} \right] = (M_n a_*, a_*) \gg \frac{1}{\lambda_n^{(1)}} \|a_*^{(n)}\|^2.$$

Отсюда $\|a_*^{(n)}\| \leq \sqrt{C}$. Теорема доказана.

Для классического метода Рунца теорема Ш.З.1 доказана в работе автора [4]; общий случай рассмотрен здесь впервые.

4°. Теорема Ш.З.2. Для того, чтобы процесс (Ш.З.3) вычисления приближенных решений по Рунцу был устойчив в последовательности пар пространств $(H_A^{(n)}, R_n)$, необходимо и достаточно, чтобы координатная система была сильно минимальна в H_A .

Процесс (Ш.З.3) сходящийся в H_A . По следствию Ш.2.1, для доказательства теоремы Ш.З.2 необходимо и достаточно установить что нормы $\|A_n^{-1}\|_{R_n \rightarrow H_A}$ от аничеши. В данном случае следует отметить оператор A_n о $M_n S_n = M_n P_n^{-1}$, так что $A_n^{-1} = P_n M_n^{-1}$.

Имеем $x_*^{(n)} = P_n a_*^{(n)}$. Как мы видели выше, $|x_*^{(n)}| = (M_n a_*, a_*)^{(n)}$

поэтому

$$|P_n a_*^{(n)}|^2 = (M_n a_*, a_*)^{(n)} = \|M_n^{-1/2} a_*^{(n)}\|^2.$$

Но $a_*^{(n)} = M_n^{-1} f^{(n)}$ и, следовательно,

$$|P_n M_n^{-1} f^{(n)}|^2 = \|M_n^{-1/2} f^{(n)}\|_{R_n}^2,$$

или

$$\|P_n M_n^{-1}\|_{R_n \rightarrow H_A^{(n)}} = \|M_n^{-1/2}\|_{R_n} = [\lambda_n^{(1)}]^{-1/2}.$$

Для того, чтобы последняя величина была ограничена, необходимо и достаточно, чтобы координатная система была сильно минимальна в энергетической норме. Теорема доказана.

Для классического метода Рунца достаточность условия теоремы Ш.З.2 доказана в статье Яковской и Яковлева [1], а необходимость - в книге автора [8]. В общем случае метод доказательства отличен от изложенного выше.

Замечание. Устойчивость вычислительных процессов, к которым приводит метод Рунца в задаче о спектре самосопряженного оператора (спектр предполагается дискретным), для классического метода Рунца исследован в статье Додина [1]. Результаты этой статьи

с некоторыми изменениями изложены в книге автора [8].

5°. Исследуем устойчивость процесса Рунда в общем случае, когда $\inf \lambda_n^{(0)} > 0$. В последовательностях пар пространств (R_n, R_n) и (H_n, R_n) эти процессы могут быть неустойчивыми (если $\inf \lambda_n^{(0)} = 0$), и мы введем вместо R_n новые пространства.

Пусть $r(n)$ — такая положительная функция от n , что $\lambda_n^{(0)} \geq c_0 r(n)$, $c_0 = \text{const}$. В частности, можно положить $r(n) = \lambda_n^{(0)}$.

Введем N -мерные гильбертовы пространства \bar{X}_n и \bar{Y}_n , элементы которых суть N -компонентные числовые векторы, а нормы заданы формулами

$$\forall v \in R_n, \|v\|_{\bar{X}_n} = r(n) \|v\|_{R_n}, \quad \|v\|_{\bar{Y}_n} = \frac{1}{r(n)} \|v\|_{R_n}. \quad (Ш.3.4)$$

Обозначим на этот раз через A_n оператор, порожденный матрицей M_n и действующий из \bar{X}_n в \bar{Y}_n , и будем считать, что $A_n^{(0)} \in \lambda_n^{(0)}$.

Теорема Ш.3.2. Процесс (Ш.3.1) устойчив в последовательности пар пространств (\bar{X}_n, \bar{Y}_n) .

Достаточно проверить, что выполнены неравенства

$$\|A_n^{-1}\|_{\bar{Y}_n \rightarrow \bar{X}_n} \leq C_1, \quad \|A_n^{(0)}\|_{\bar{X}_n} \leq C_2. \quad (Л.3.5)$$

Оценим величину $\|A_n^{-1}\|_{\bar{Y}_n \rightarrow \bar{X}_n}$:

$$\begin{aligned} \|A_n^{-1}\|_{\bar{Y}_n \rightarrow \bar{X}_n} &= \\ &= \sup_{b \in R_n} \frac{\|A_n^{-1} b\|_{\bar{X}_n}}{\|b\|_{\bar{Y}_n}} = r^2(n) \sup_{b \in R_n} \frac{\|M_n^{-1} b\|_{R_n}}{\|b\|_{R_n}} = \frac{r^2(n)}{\lambda_n^{(0)}} \leq \frac{1}{c_0^2}, \end{aligned}$$

и первое из неравенств (Л.3.5) устан. влнво. Заметим, что если принять $r(n) = \sqrt{\lambda_n^{(0)}}$, то получится $\|A_n^{-1}\|_{\bar{Y}_n \rightarrow \bar{X}_n} = 1$.

Пусть $v^{(n)}$ — произвольный элемент из $H_n^{(n)}$:

$$v^{(n)} = \sum_{k=1}^N b_k y_k \dots$$

полагаем $b = (b_1, b_2, \dots, b_N)$, имеем

$$\|v^{(n)}\|^2 = (M_n b, b)_{R_n} \geq \lambda_n^{(0)} \|b\|_{R_n}^2 \geq c_0^2 \|b\|_{R_n}^2. \quad (П.3.7)$$

Метод Рунда сходится в H_n , поэтому $\|A_n^{(0)}\| \leq C = \text{const}$.

По неравенству (Ш.3.6) $\|a_*^{(n)}\| \leq C/c_0$ и второе неравенство (Ш.3.5) также установлено. Теорема доказана.

Замечание Ш.3.1. Пусть $\mu_n^{(j)}$ - наименьшее собственное число матрицы \bar{M}_n скалярных произведений $(y_{nk}, y_{nj})_H$; $j, k = 1, 2, \dots, N$. Достаточно подчинить $\gamma^{(n)}$ неравенству $\mu_n^{(j)} \geq c^2 \gamma^2(n)$, $c = \text{const} > 0$. Для доказательства достаточно заметить, что $\mu_n^{(j)} \leq c' \lambda_n^{(j)}$, $c' = \text{const}$. Действительно,

$$\mu_n^{(j)} = \inf_{b \in R_{jX}} \frac{(\bar{M}_n b, b)_{R_{jX}}}{\|b\|_{R_{jX}}^2} = \inf_{b \in R_{jX}} \frac{\left\| \sum_{k=1}^N \bar{b}_k y_{nk} \right\|_H^2}{\|b\|_{R_{jX}}^2},$$

$$\lambda_n^{(j)} = \inf_{b \in R_{jX}} \frac{(M_n b, b)_{R_{jX}}}{\|b\|_{R_{jX}}^2} = \inf_{b \in R_{jX}} \frac{\left| \sum_{k=1}^N b_k y_{nk} \right|^2}{\|b\|_{R_{jX}}^2}.$$

Пусть λ_0 - нижняя грань спектра оператора A , и пусть второй инфимум достигается при $b = \bar{b} = (b_1, b_2, \dots, b_N)$.

Тогда

$$\mu_n^{(j)} \leq \frac{\left\| \sum_{k=1}^N \bar{b}_k y_{nk} \right\|_H^2}{\|\bar{b}\|_{R_{jX}}^2} \leq \frac{1}{\lambda_0} \frac{\left| \sum_{k=1}^N \bar{b}_k y_{nk} \right|^2}{\|\bar{b}\|_{R_{jX}}^2} = \frac{\lambda_n^{(j)}}{\lambda_0},$$

и можно положить $c' = 1/\lambda_0$.

6°. Теорема Ш.3.4. Процесс (Ш.3.3) устойчив в последовательности пар пространств $(H_A^{(n)}, \bar{Y}_n)$.

Для доказательства этой теоремы мы воспользуемся приемом статьи Локковой и Яковлева [1].

Уравнение (Ш.3.1) запишем в виде

$$A_n a^{(n)} = f^{(n)}, \quad (\text{Ш.3.7})$$

где мы считаем, что $a^{(n)} \in \bar{X}_n$, $f^{(n)} \in \bar{Y}_n$, и обозначаем через A_n оператор, действующий из \bar{X}_n в \bar{Y}_n и порожденный матрицей M_n . Наряду с уравнением (Ш.3.7) рассмотрим искаженное уравнение

$$(A_n + \Gamma_n) C^{(n)} = f^{(n)} + \delta^{(n)}. \quad (\text{Ш.3.8})$$

Исканное приближенное решение по Рунту имеет вид

$$\bar{x}_*^{(n)} = \sum_{k=1}^N C_k^{(n)} y_{nk} = P_n C^{(n)}.$$

Следовательно

$$\begin{aligned} \|z_*^{(n)} - x_*^{(n)}\| &= (M_n(c^{(n)} - a^{(n)}), c^{(n)} - a^{(n)})_{R_{\mathcal{H}}} \leq \\ &\leq \|A_n(c^{(n)} - a^{(n)})\|_{\bar{y}_n} \cdot \|c^{(n)} - a^{(n)}\|_{\bar{x}_n}. \end{aligned} \quad (\text{Ш.3.9})$$

В силу теоремы Ш.3.3 существуют такие числа $\rho, q, \nu > 0$, что если $\|\Gamma_n\|_{\bar{x}_n \rightarrow \bar{y}_n} \leq \nu$, то

$$\|c^{(n)} - a^{(n)}\|_{\bar{x}_n} \leq \rho \|\Gamma_n\|_{\bar{x}_n \rightarrow \bar{y}_n} + q \|\delta^{(n)}\|_{\bar{y}_n}. \quad (\text{Ш.3.10})$$

Из (Ш.3.7) и (Ш.3.8) следует

$$(A_n + \Gamma_n)(c^{(n)} - a^{(n)}) = (I_n + \Gamma_n A_n^{-1}) A_n (c^{(n)} - a^{(n)}) = \delta^{(n)} - \Gamma_n a^{(n)},$$

I_n - тождественный оператор в \bar{y}_n . Как было доказано в п.5⁰, $\|A_n^{-1}\|_{\bar{y}_n \rightarrow \bar{x}_n} \leq 1/c_0^2$. Потребуем, чтобы $\|\Gamma_n\|_{\bar{x}_n \rightarrow \bar{y}_n} \leq \beta/c_0^2$, где β - какое-нибудь число из интервала $(0, 1)$.

Тогда $\|(I_n + \Gamma_n A_n^{-1})\|_{\bar{y}_n} \leq (1-\beta)^{-1}$ и, так как нормы $\|a^{(n)}\|_{\bar{x}_n}$ ограничены, то

$$\|A_n(c^{(n)} - a^{(n)})\|_{\bar{y}_n} \leq \frac{1}{1-\beta} [\|\delta^{(n)}\|_{\bar{y}_n} + C' \|\Gamma_n\|_{\bar{x}_n \rightarrow \bar{y}_n}], \quad C' = \text{const}. \quad (\text{Ш.3.11})$$

Из формул (Ш.3.9) - (Ш.3.11) очевидно следует, что

$$\|z_*^{(n)} - x_*^{(n)}\| \leq \rho_1 \|\Gamma_n\|_{\bar{x}_n \rightarrow \bar{y}_n} + q_1 \|\delta^{(n)}\|_{\bar{y}_n}; \quad (\text{Ш.3.12})$$

за ν_1 можно принять число $\nu_1 = \min(\nu, \beta/c_0^2)$.

7⁰. Частными случаями теорем Ш.3.3 и Ш.3.4 являются сформулированные ниже, в гл.V, теоремы об устойчивости процесса конечных элементов.

8⁰. Пусть $\inf_n \lambda_n^{(0)} = 0$. Тогда процесс вычисления коэффициентов Рунда неустойчив в $(R_{\mathcal{H}}, R_{\mathcal{H}})$, а процесс вычисления приближенного решения неустойчив в $(H_A^{(n)}, R_{\mathcal{H}})$. Это значит, что погрешность $\|c^{(n)} - a^{(n)}\|_{R_{\mathcal{H}}}$ и $\|z_*^{(n)} - x_*^{(n)}\|$ могут быть велики, даже если малы погрешности $\|\Gamma_n\|_{R_{\mathcal{H}} \rightarrow R_{\mathcal{H}}}$ и $\|\delta^{(n)}\|_{R_{\mathcal{H}}}$. Нетрудно, однако, убедиться, что оценка для $\|z_*^{(n)} - x_*^{(n)}\|$ в этом случае оказывается существенно меньшей, чем для величины $\|c^{(n)} - a^{(n)}\|_{R_{\mathcal{H}}}$ действитель-

но, из (Ш.3.10) и (Ш.3.12) следует

$$\begin{aligned} \|c^{(n)} - a^{(n)}\|_{R_f} &\leq \frac{1}{r(n)} \left[\frac{p}{r^2(n)} \|\Gamma_n\|_{R_f \rightarrow R_f} + \frac{q}{r(n)} \|\delta^{(n)}\|_{R_f} \right] \leq \\ &\leq \frac{\bar{p}}{r(n)} \left[\frac{\|\Gamma_n\|_{R_f \rightarrow R_f}}{r^2(n)} + \frac{\|\delta^{(n)}\|_{R_f}}{r(n)} \right]; \quad \bar{p} = \max(p, q); \\ \|z_n^{(n)} - c_n^{(n)}\| &\leq \left[\frac{p'}{r^2(n)} \|\Gamma_n\|_{R_f \rightarrow R_f} + \frac{q'}{r(n)} \|\delta^{(n)}\|_{R_f} \right] \leq \\ &\leq \bar{p}' \left[\frac{\|\Gamma_n\|_{R_f \rightarrow R_f}}{r^2(n)} + \frac{\|\delta^{(n)}\|_{R_f}}{r(n)} \right]; \quad \bar{p}' = \max(p', q'); \end{aligned}$$

правая часть первого неравенства содержит малый делитель $r(n)$, откуда и следует наше утверждение.

§ 4. Примеры.

1°. Пример 1. Рассмотрим уравнение

$$(t+1)x(t) = 1, \quad 0 \leq t \leq 1. \quad (\text{Ш.4.1})$$

Оператор умножения на $t+1$ - положительно определенный в $L_2(0,1)$, и уравнение (Ш.4.1) равносильно задаче с экстремуме функционала

$$\int_0^1 [(t+1)x^2(t) - 2x(t)] dt. \quad (\text{Ш.4.2})$$

Эту последнюю задачу можно решать по методу Рунца. В качестве координатной системы возьмем последовательность t^n , $n=0,1,2,\dots$. По известной теореме Мюнца (см., например, Натансон [1]) эта система полна в $L_2(0,1)$, так как ряд обратных показателей $\sum_{n=1}^{\infty} \frac{1}{n}$ расходится. Из той же теоремы Мюнца легко вытекает (см. книгу автора [8]), что эта же система не сильно минимальна и процесс Рунца с такой координатной системой неустойчив в $L_2(0,1)$. Посмотрим, как эта неустойчивость проявится в счете.

2°. При выбранной выше координатной системе приближенное решение имеет вид

$$x^{(n)}(t) = \sum_{k=0}^n a_k^{(n)} t^k, \quad (\text{Ш.4.3})$$

что приводит к неосвоенной системе Рунца

$$\sum_{j=0}^n \frac{2k-2j+3}{(k+j+1)(k+j+2)} a_j^{(n)} = \frac{1}{j+1}, \quad 0 \leq j \leq n-1. \quad (\text{Ш.4.4})$$

Напишем расширенную матрицу Рунца 5-го порядка:

<u>3</u>	<u>5</u>	<u>7</u>	<u>9</u>	<u>11</u>	<u>1</u>
2	6	12	20	30	

<u>5</u>	<u>7</u>	<u>9</u>	<u>11</u>	<u>13</u>	<u>1</u>
6	12	20	30	42	2

<u>7</u>	<u>9</u>	<u>11</u>	<u>13</u>	<u>15</u>	<u>1</u>
12	20	30	42	56	3

(П.4.5.)

<u>9</u>	<u>11</u>	<u>13</u>	<u>15</u>	<u>17</u>	<u>1</u>
20	30	42	56	72	4

<u>11</u>	<u>13</u>	<u>15</u>	<u>17</u>	<u>19</u>	<u>1</u>
30	42	56	72	90	5

Как обычно, расширенные матрицы Рунца низшего порядка получают усечением.

Системы Рунца порядка от 3 до 5 решались на ЭВМ БЭСМ-6. При этом оказалось, что для системы 5-го порядка наступает перегибание, и получить решение не удалось. Тем более не имело смысла решать системы Рунца более высоких порядков. Точные значения коэффициентов $A_k^{(n)}$, $k = 1, 2, 3, 4$ приведены в таблице 1. В таблице 2 приведены десятичные приближения тех же коэффициентов с шестью верными знаками. Приведены также десятичные приближения коэффициентов для $n = 5$, полученные следующим образом. Так как получить точные значения этих коэффициентов не удалось, то для $n = 5$ элементы расширенной матрицы Рунца были заменены их десятичными приближениями с 12 знаками; после этого были получены десятичные приближения коэффициентов, в которых мы сохранили 6 знаков.

Представляется интересным выяснить, насколько искаженное приближение близко (в энергетической матрице) к точному решению $x_4(t)$. В нашем примере энергетическая норма определяется формулой

$$\|x\|^2 = \int_0^1 (1+t)^2 x^2(t) dt$$

(П.4.6)

и эквивалентна L_2 -норме. Были вычислены значения $\|x_n - x_4^{(n)}\|$ для $n = 1, 2, 3, 4$. Результаты приведены в таблице 3.

Таблица 1. Точные значения неискаженных коэффициентов Рунца

$n \backslash k$	1	2	3	4
1	2/3			
2	12/13	- 6/13		
3	62/63	- 50/63	20/63	
4	320/321	- 100/107	70/107	- 70/321

Таблица 2. Десятичные приближения неискаженных коэффициентов Рунца

	1	2	3	4	5
1	0,333333				
2	0,923077	-0,461538			
3	0,981270	-0,793651	0,317460		
4	0,996885	-0,334579	0,654206	-0,218069	
5	0,992859	-0,864970	0,367892	0,200555	-0,199993

Таблица 3. Погрешность аппроксимации приближенных решений

n	1	2	3	4	5
	0,1527	0,0290	0,0051	0,0008	0,0016

Таблица 3 показывает, что приближенные решения, вычисляемые без погрешности довольно хорошо сходятся к точному решению в энергетической метрике. В классическом методе Рунца погрешность аппроксимации в энергетической метрике убывает с ростом n . Последнее число в таблице 3 показывает, что при $n = 5$ коэффициенты Рунца оказались вычисленными недостаточно точно.

В заключение настоящего пункта приведем таблицу значений в комод функции точных и приближенных, соответствующих значениям $n = 3, 4, 5$. Все значения даны с шестью верными знаками после запятой.

Таблица 4. Точные и приближенные неискаженные значения искомой функции

t	точные значения	приближенные значения		
		$n = 3$	$n = 4$	$n = 5$
0,0	1,000000	0,984310	0,995867	0,994871
0,1	0,909091	0,908003	0,909744	0,910432
0,2	0,833333	0,838072	0,835098	0,833902
0,3	0,769231	0,774517	0,773706	0,776198
0,4	0,714286	0,717339	0,721801	0,708751
0,5	0,666667	0,666538	0,678546	0,663192
0,6	0,625000	0,622114	0,641773	0,630532
0,7	0,588235	0,584066	0,608823	0,609512
0,8	0,555556	0,552595	0,576560	0,593628
0,9	0,526316	0,527100	0,541364	0,566283
1,0	0,500000	0,508182	0,499133	0,493396

Как мы видим, приближенные значения решений в фиксированных точках приближаются к точным значениям довольно медленно и не монотонно. Это связано с тем, что метод Рунца охотится, вообще говоря, только в энергетической метрике, которая в данном случае эквивалентна метрике L_2 .

3°. Покажем расширенную матрицу y Рунца, заменив все ее элементы их десятичными приближениями с четырьмя верными знаками после запятой. Приведем таблицу искаженных коэффициентов Рунца, вычисленных с шестью верными знаками после запятой.

Таблица 5. Искаженные коэффициенты Рунца

$n \backslash k$	1	2	3	4	5
3	0,984310	-0,794960	0,517048		
4	0,995867	-0,923202	0,627642	-0,201115	
5	0,994871	-0,852381	0,377333	0,258677	-0,255133

Сравнивая эт. с таблицей 2, видим, что искаженные коэффициенты оказываются уже, как правило, во втором или третьем знаке после запятой, тогда как в расширенной матрице Рунца и как-то

возмущение только пятым знаком. Это связано с неустойчивостью процесса Рунге в данном примере в последовательности пространств (R_{n-1}, R_n) .

Рассмотрим теперь процесс решения задачи (Ш.4.1) по Рунге, когда за координатные функции взяты полиномы Лежандра, умноженные на множители, по порядку роста совпадающие с нормируемыми множителями:

$$y_n(t) = [\sqrt{n}] P_{n-1}(2t-1).$$

Функции $y_n(t)$ сильно минимальны в $L_2(0,1)$ (см. книгу автора [8]) а, следовательно, и в энергетической метрике задачи (Ш.4.1). При $n \leq 6$ получены точные (в виде простых дробей) значения Рунге; в таблице 6 приведены десятичные значения этих коэффициентов с шестью верными знаками после запятой. При взгляде на эту таблицу можно догадаться, что коэффициенты Рунге $a_k^{(n)}$ при возрастании n стремятся к некоторым пределам, а при возрастании k довольно быстро убывают. В книге автора [5] показано, что если координатная система классического метода Рунге сильно минимальна в энергетической норме, то осуществуют пределы $a_k = \lim_{n \rightarrow \infty} a_k^{(n)}$ и последовательность $\{a_k\} \in L_2$.

Расширенные матрицы были подвергнуты искажению, такому же, как в п. 3^о: их элементы были заменены десятичными приближениями с четырьмя первыми знаками, и полученные системы были решены с возможно большей точностью. Значения искаженных коэффициентов Рунге с шестью верными знаками приведены в таблице 7. Сравнение с таблицей 6 показывает, что искажения коэффициентов Рунге проявляются только в пятом и шестом знаках и, следовательно, являются величинами того же порядка, что и искажения элементов расширенной матрицы Рунге.

Таблица 6. Десятичные приближения немасштабированных коэффициентов Рунге; координатные функции - полиномы Лежандра

	1	2	3	4	5	6
5	0,699116	-0,238918	0,054516	-0,005432		
1	0,699147	-0,238924	0,054565	-0,005615	0,001073	
2	0,699147	-0,236325	0,054567	-0,005620	0,001101	-0,000204
3	0,699147	-0,236625	0,054567	-0,005620	0,001102	-0,000210

Значения $a_6^{(5)} = 0,000038$

Таблица 7. Десятичные приближения искаженных коэффициентов Ритца; координатные функции - полиномы Лагранжа

	1	2	3	4	5	6
3	0,693146	-0,238314	0,054516	-0,005452		
4	0,693146	-0,238320	0,054564	-0,005615	0,001070	
5	0,693146	-0,238321	0,054566	-0,005620	0,001101	-0,000204
6	0,693146	-0,238321	0,054566	-0,005620	0,001102	-0,000210

Искаженные значения $a_3^{(n)} = 0,000038$.

5°. В пп. 5° и 6° настоящего параграфа мы кратко опишем результаты вычислений по двум примерам, подробно рассмотренным в книге автора [8]. Будем решать по методу Ритца задачу, положительно определенную в $L_2(0, 1)$:

$$0 \leq t \leq 1, \quad -\frac{d^2 x}{dt^2} = \frac{1}{1+t}; \quad x(0) = x'(1) = 0; \quad (\text{Ш.4.6})$$

В качестве координатной системы возьмем последовательность $\varphi_n(t) = t^n$, $n = 1, 2, \dots$; опираясь на теорему Монца, легко доказать, что эта система полна, но не сильно минимальна в энергетической метрике задачи (Ш.4.6).

Первые восемь приближений можно получить, решая линейную алгебраическую систему с расширенной матрицей

I	I	I	I	I	I	I	I	I	I	-	ln 2
I	4/3	3/2	8/5	5/3	12/7	7/4	16/9		ln 2	-	1/2
I	3/2	9/5	2	15/7	9/4	7/3	12/5		5/6	-	ln 2
I	8/5	2	16/7	5/2	8/3	14/5	32/11		ln 2	-	7/12
I	5/3	15/7	5/2	25/9	3	35/11	10/3		47/6	-	ln 2
I	12/7	9/4	8/3	3	36/11	7/2	48/13		ln 2	-	37/60
I	7/4	7/3	14/5	35/11	7/2	49/13	4		319/420	-	ln 2
I	16/9	12/5	32/11	10/3	48/13	4	64/5		ln 2	-	539/840

а также все системы, полученные из этой усечением. Коэффициенты Ритца - решения упомянутых систем - обозначим через $a_k^{(n)}$. Решим эти точно, без погрешностей, а затем за ним полученные таким образом искаженные значения коэффициентов Ритца их десятичными приближениями с четырьмя верными знаками (таблица 8).

Таблица 8. Двоичные приближения неискаженных коэффициентов Ритца; координатные функции - степени t .

	1	2	3	4	5	6	7	8
1	0,3069							
2	0,6480	-0,3411						
3	0,6869	-0,4579	0,0788					
4	0,6922	-0,4899	0,1312	-0,0267				
5	0,6930	-0,4977	0,1547	-0,0541	0,0110			
6	0,6931	-0,4995	0,1631	-0,0708	0,0260	-0,0050		
7	0,6931	-0,4999	0,1657	-0,0786	0,0378	-0,0136	0,0025	
8	0,6931	-0,5000	0,1664	0,0817	0,0446	0,0218	0,0075	0,0013

Искажем теперь расширенную матрицу (Ш.4.7), заменив в ней свободные члены их двоичными приближениями с четырьмя верными знаками после запятой. Искаженные коэффициенты Ритца, которые получим, решая систему с искаженной расширенной матрицей, а также соответствующие усеченные системы, обозначим через $b_k^{(n)}$. Но системы опять решим точно, а затем заменим коэффициенты $b_k^{(n)}$ двоичными приближениями с четырьмя верными знаками (таблица

Таблица 9. Десятичные приближения искаженных коэффициентов Ритца

	1	2	3	4	5	6	7	8
1	0,3069							
2	0,6483	-0,3414						
3	0,6883	-0,4614	0,0800					
4	0,6974	-0,5160	0,1710	-0,0455				
5	0,7127	-0,6690	0,6300	-0,5810	0,2142			
6	0,7635	-1,128	4,172	-7,665	6,590	-2,125		
7	0,9141	-4,595	25,28	-71,00	101,6	-71,79	19,91	
8	1,627	-24,54	204,8	-819,1	1747	-2047	1243	-305,7

Сравнение таблиц 8 и 9 показывает, что погрешность искажения коэффициентов Ритца в данном случае весьма сильно возрастает вместе с n .

6°. Рассмотрим еще задачу

$$-1 \leq t \leq 1, \quad -\frac{d}{dt} \left[(2+t) \frac{dx}{dt} \right] = 1, \quad x(-1) = x(1) = 0,$$

положительно определенную в $L_2(0, 1)$. В качестве координатной системы выберем последовательность интегралов от нормированных полиномов Лагранжа

$$y_n(t) = \sqrt{\frac{2n+1}{2}} \int_{-1}^t P_n(t) dt, \quad n=1, 2, \dots \quad (Ш.4.9)$$

Эта система сильно минимальна в энергетическом пространстве задачи (Ш.4.8) (см. книгу автора [8]). Как и в предшествующем примере, были составлены и точно решены системы Рунца от 1-го до n -го порядка. В таблице IО приведены десятичные приближения этих решений с четырьмя верными знаками после запятой. Далее системы Рунца были искажены: все элементы расширенных матриц были заменены их десятичными приближениями с четырьмя верными знаками после запятой. В этом случае получение точных решений оказалось затруднительным, поэтому, чтобы по возможности избежать ошибок округления, вычисления производились с большим числом десятичных знаков - с семью знаками для систем до 4-го порядка включительно и с 12 знаками для систем порядка от 5-го по 8-й.

В таблице II приведены искаженные значения коэффициентов Рунца с четырьмя верными знаками после запятой. Сравнение двух последних таблиц показывает, что в данном случае искажения коэффициентов Рунца суть величины того же порядка, что и искаженные элементы расширенной матрицы Рунца.

Таблица IО. Десятичные приближения неискаженных коэффициентов Рунца.

	1	2	3	4	5	6	7	8
1	-0,4082							
2	-0,4374	0,1129						
3	-0,4395	0,1213	-0,0008					
4	-0,4397	0,1219	-0,0030	0,0083				
5	-0,4397	0,1219	-0,0032	0,0089	-0,0022			
6	-0,4397	0,1220	-0,0032	0,0090	-0,0024	0,0006		
7	-0,4397	0,1220	-0,0032	0,0090	-0,0024	0,0006	-0,0002	
8	-0,4397	0,1220	-0,0032	0,0090	-0,0024	0,0006	-0,0002	0,0001

70. См. один пример, в котором отмечаются явления чувствительности, дан в работе автора [17] и [2.].

Таблица II. Десятичные приближения искаженных коэффициентов Ритца.

	1	2	3	4	5	6	7	8
1	-0,4083							
2	-0,4374	0,1129						
3	-0,4396	0,1213	-0,0308					
4	-0,4397	0,1219	-0,0330	0,0083				
5	-0,4398	0,1220	-0,0332	0,0089	-0,0022			
6	-0,4398	0,1220	-0,0332	0,0089	-0,0024	0,0006		
7	-0,4398	0,1220	-0,0332	0,0089	-0,0024	0,0007	-0,0002	
8	-0,4398	0,1220	-0,0332	0,0089	-0,0024	0,0007	-0,0002	0,0001

§ 5.06 устойчивости процесса Бубнова - Галеркина

1°. В сепарабельном гильбертовом пространстве, которое для упрощения обозначений считаем вещественным, рассмотрим уравнение

$$Ax = A_0x + Kx = f, \quad (Ш.5.1)$$

где A_0 - положительно определенный самоопределенный оператор и произведение $T = A_0^{-1}K$ вполне непрерывно в энергетическом пространстве H_0 оператора A_0 . В работе автора 1948 г. (подробнее см. книгу автора [3]) установлены следующие теоремы:

а) Уравнение (3.5.1) эквивалентно уравнению

$$x + Tx = \tilde{A}^{-1}f = \tilde{f}, \quad (Ш.5.2)$$

рассматриваемому в H_0 , в следующем смысле: любое решение уравнения (Ш.5.1) удовлетворяет уравнению (Ш.5.2), а любое решение уравнения (Ш.5.2) является обобщенным решением уравнения (Ш.5.1).

б) При заданной координатной системе метод Бубнова - Галеркина, примененный к уравнениям (Ш.5.1) и (Ш.5.2), приводит к одной и той же алгебраической системе.

в) Если координатная система полна в H_0 , то метод Бубнова - Галеркина для уравнения (Ш.5.1) (а, следовательно, и для уравнения (Ш.5.2)) сходится в H_0 .

2°. Из сказанного в п.1° следует, что устойчивость метода Бубнова - Галеркина достаточно исследовать для уравнения (Ш.5.1). Как и в случае метода Ритца, мы рассмотрим два вычислительных процесса: для вычисления вектора коэффициентов Бубнова - Галер-

кина и для вычисления приближенного решения по Бубнову - Галеркину. Устойчивость этих процессов исследована в работе Яковой и Дювлева [1] для более общего метода, часто называемого методом Петрова - Галеркина; доказано, что для устойчивости названных процессов достаточно, чтобы координатная система была сильно минимальна в H_0 . В статье Вайникко [7] доказано, что условие сильной минимальности и необходимо; там же упрощено доказательство достаточности. В данном параграфе мы изложим содержание статьи Вайникко [7], ограничившись при этом более простым методом Бубнова - Галеркина.

3°. Пусть для построения приближенного по Бубнову - Галеркину решения уравнения (Ш.5.2) использована координатная система

$\{y_n\} \in H_0$, подчиненная обычным условиям; в частности, при любом n элементы y_1, y_2, \dots, y_n предполагаются линейно независимыми. Систему $\{y_n\}$ можно ортонормировать в H_0 ; пусть $\{\omega_n\}$ - система, полученная ортогонализацией системы $\{y_n\}$. Обе последовательности связаны рекуррентными соотношениями вида

$$\omega_k = \sum_{j=1}^k c_{kj} y_j, \quad y_k = \sum_{j=1}^k \gamma_{kj} \omega_j. \quad (\text{Ш.5.3})$$

Обозначим через C_n матрицу первых n соотношений (Ш.5.3).

Лемма Ш.5.1. Если последовательность $\{y_n\}$ сильно минимальна в H_0 , то нормы $\|C_n\|_{R_n \rightarrow R_n}$ ограничены; в противном случае они монотонно стремятся к бесконечности вместе с n .

Имеем $[\omega_k, \omega_j] = \delta_{kj}$; используя представление (Ш.5.3), получаем

$$C_n M_n C_n^* = I_n, \quad (\text{Ш.5.4})$$

где M_n имеет тот же смысл, что и в § 3, звездочка означает сопряженную матрицу. I_n - единичная матрица в R_n . Из (Ш.5.4) следует, что $M_n^{-1} = C_n^* C_n$ и, следовательно, $\|C_n\| = \sqrt{\|M_n^{-1}\|} = \|\lambda_n^{(1)}\|^{-1/2}$. Дальнейшее очевидно.

Теорема Ш.5.1. Для устойчивости процесса вычисления приближенного по Бубнову - Галеркину решения уравнения (Ш.5.2) в последовательности пар пространств $(H_0^{(n)}, R_n)$ необходимо и достаточно, чтобы координатная система была сильно минимальна в H_0 .

Здесь $H_0^{(n)}$ обозначено подпространство H_0 с базисом (y_i) .

y_2, \dots, y_n).

Обозначим соответственно неискаженный и искаженный вектор коэффициентов Бубнова - Галеркина через $a_*^{(n)}$ и $b_*^{(n)}$, неискаженное и искаженное приближенное решение - через $x_*^{(n)}$ и $\tilde{x}_*^{(n)}$. Далее, положим

$$f^{(n)} = ((f, y_1), (f, y_2), \dots, (f, y_n)), y^{(n)} = (y_1, y_2, \dots, y_n), d^{(n)} = (C_n^*)^{-1} a_*^{(n)}, \beta^{(n)} = (C_n^*)^{-1} b_*^{(n)}.$$

Наконец, через B_n обозначим матрицу элементов

$$[y_k + T y_k, y_j]; \quad j, k = \overline{1, 2, \dots, n}. \text{ В этих обозначениях имеем}$$

$$\text{Далее, } a_*^{(n)} = C_n^* d^{(n)}, \text{ откуда}$$

$$x_*^{(n)} = \sum_{k=1}^n a_k^{(n)} y_k = a_*^{(n)} \cdot y^{(n)}.$$

$$x_*^{(n)} = \sum_{k=1}^n \sum_{j=k}^n C_{jk} d_j^{(n)} y_k = \sum_{j=1}^n d_j^{(n)} \sum_{k=1}^j C_{jk} y_k = \sum_{j=1}^n d_j^{(n)} \omega_j.$$

Аналогично

$$\tilde{x}_*^{(n)} = \sum_{j=1}^n \beta_j^{(n)} \omega_j$$

и, следовательно,

$$\tilde{x}_*^{(n)} - x_*^{(n)} = \sum_{j=1}^n (\beta_j^{(n)} - d_j^{(n)}) \omega_j. \quad (\text{Ш.5.5})$$

Системы уравнений Бубнова - Галеркина, искаженная и неискаженная, имеют соответственно вид

$$(B_n + \Gamma_n) b^{(n)} = f^{(n)} + \delta^{(n)}; \quad B_n a^{(n)} = f^{(n)},$$

или

$$(C_n B_n C_n^* + C_n \Gamma_n C_n^*) \beta^{(n)} = C_n (f^{(n)} + \delta^{(n)}),$$

$$C_n B_n C_n^* d^{(n)} = C_n f^{(n)}.$$

Отсюда

$$(C_n B_n C_n^* + C_n \Gamma_n C_n^*) (\beta^{(n)} - d^{(n)}) = C_n \delta^{(n)} - C_n \Gamma_n C_n^* d^{(n)}. \quad (\text{Ш.5.6})$$

Обозначим через Π_n ортопроектор из H_0 на $H_0^{(n)}$. Докажем, что существует такая постоянная $\tau > 0$, что

$$\forall n, \forall v^{(n)} \in (I+T)H_0^{(n)}; |\Pi_n v^{(n)}| \geq \tau |v^{(n)}|. \quad (II.5.7)$$

Допустим противное. Тогда можно построить такую последовательность $v^{(n)} \in (I+T)H_0^{(n)}$, $n = 1, 2, \dots$, что $|v^{(n)}| = 1$ и $|\Pi_n v^{(n)}| \xrightarrow{n \rightarrow \infty} 0$. Пусть $v^{(n)} = (I+T)w^{(n)}$, тогда $w^{(n)} \in H_0^{(n)}$. При этом $|w^{(n)}| \leq |I+T|$ и $|\Pi_n (I+T)w^{(n)}| = |w^{(n)} + \Pi_n T w^{(n)}| \xrightarrow{n \rightarrow \infty} 0$.

Последовательность $\{w^{(n)}\}$ ограничена и потому слабо компактна в H_0 . Выделим из $\{w^{(n)}\}$ подпоследовательность, которую опять обозначим через $\{w^{(n)}\}$ и которая слабо сходится в H_0 к некоторому элементу $w^{(0)}$. Тогда $T w^{(n)} \rightarrow T w^{(0)}$ в H_0 . Далее

$$\Pi_n T w^{(n)} = \Pi_n T w^{(0)} + \Pi_n T (w^{(n)} - w^{(0)}).$$

Первое слагаемое справа стремится к $T w^{(0)}$ в норме H_0 , потому что последовательность подпространств $\{H_0^{(n)}\}$ полна в H_0 ; второе слагаемое стремится к нулю в той же норме. Таким образом, $\Pi_n T w^{(n)} \rightarrow T w^{(0)}$. Теперь из соотношения $|w^{(n)} + \Pi_n T w^{(n)}| \rightarrow 0$ следует, что $(I+T)w^{(0)} = 0$. Отсюда вытекает, что

$$\lim v^{(n)} = \lim (I+T)w^{(n)} = (I+T)w^{(0)} = 0,$$

что противоречит условию $|v^{(n)}| = 1$.

Докажем теперь, что при достаточно больших n имеет место неравенство

$$\begin{aligned} \forall c^{(n)} \in R_n, \frac{1}{\mu} \|c^{(n)}\| &\leq \|C_n B_n C_n^* c^{(n)}\| \leq \\ &\leq \|I+T\| \cdot \|c^{(n)}\|, \quad \mu = \tau^{-1} |(I+T)^{-1}|. \end{aligned} \quad (II.5.8)$$

Как легко видеть, $C_n B_n C_n^*$ есть матрица скалярных произведений $[\omega_k + T\omega_k, \omega_j]$; $j, k = 1, 2, \dots, n$. Отсюда

$$\begin{aligned} \|C_n B_n C_n^* c^{(n)}\|^2 &= \sum_{j=1}^n \left| \sum_{k=1}^n (\omega_k + T\omega_k, \omega_j) c_k^{(n)} \right|^2 = \\ &= \sum_{j=1}^n \left| \left((I+T) \sum_{k=1}^n c_k^{(n)} \omega_k, \omega_j \right) \right|^2 = \left| \Pi_n (I+T) \sum_{k=1}^n c_k^{(n)} \omega_k \right|^2 \leq \end{aligned}$$

$$\leq \|I+T\|^2 \cdot \left\| \sum_{k=1}^n c_k^{(n)} \omega_k \right\|^2 = \|I+T\|^2 \cdot \|c^{(n)}\|.$$

С другой стороны, неравенство (Ш.5.7) дает

$$\mu_n (I+T) \sum_{k=1}^n c_k^{(n)} \omega_k \geq \tau^2 (I+T) \sum_{k=1}^n c_k^{(n)} \omega_k \geq \frac{1}{\mu^2} \left\| \sum_{k=1}^n c_k^{(n)} \omega_k \right\|^2 = \frac{1}{\mu^2} \|c^{(n)}\|^2.$$

4°. Переходим к доказательству утверждений теоремы. Пусть координатная система $\{y_n\}$ сильно минимальна в H_n . По лемме Ш.5.1 $\|C_n\| = \|C_n^*\| \leq c_0 = \text{const}$. Положим $\nu = (2\mu c_0^2)^{-1}$, и пусть $\|\Gamma_n\| \leq \nu$. Тогда $\|C_n \Gamma_n C_n^*\| \leq 1/2\mu$. В то же время, как показывает левое неравенство (Ш.5.8), $\|C_n B_n C_n^*\| \geq 1/\mu$,

матрица в левой части уравнения (Ш.5.6) обратима. При этом $\| [C_n (B_n + \Gamma_n) C_n^*]^{-1} \| \leq 2\mu$; из соотношений (Ш.5.6) и из факта сходимости метода Бубнова - Галеркина следует

$$\|x_n^{(n)} - x_n^*\| = \|\beta^{(n)} - a^{(n)}\| \leq \rho \|\Gamma_n\| + q \|\delta^{(n)}\|, \quad (\text{Ш.5.9})$$

где ρ и q , некоторые постоянные. Достаточность условия теоремы доказана.

Пусть теперь координатная система не сильно минимальна в H_n . По лемме Ш.5.1 $\|C_n\| \xrightarrow{n \rightarrow \infty} \infty$. Выберем $\delta^{(n)}$ так, чтобы

$$\|C_n \delta^{(n)}\| \geq \frac{1}{2} \|C_n\| \cdot \|\delta^{(n)}\|. \quad \text{Правое неравенство (Ш.5.3) показывает, что при } \Gamma_n = 0$$

$$\|\beta^{(n)} - a^{(n)}\| \geq \frac{\|C_n\|}{2\|I+T\|} \|\delta^{(n)}\|$$

и, так как коэффициент при $\|\delta^{(n)}\|$ бесконечно возрастает, то рассматриваемый процесс неустойчив. Этим доказана необходимость условия теоремы.

5°. Теорема Ш.5.2. Если координатная система сильно минимальна в H_n , то процесс вычисления коэффициентов Бубнова - Галеркина устойчив в последовательности пар пространств (R_n, R_n) .

Доказательство очень просто: так как $b^{(n)} - a^{(n)} = C_n (\beta^{(n)} - a^{(n)})$, $\|C_n\| \leq c_0$, то по неравенству (Ш.5.9)

$$\|b^{(n)} - a^{(n)}\| \leq c_0 \rho \|\Gamma_n\| + c_0 q \|\delta^{(n)}\|.$$

Замечание. В работах М.А. Велиева [1 - 9] после года устоячивость метода Бубнова - Галеркина (точнее, Фаздо - Галеркина) для довольно широкого класса нестационарных операторных уравнений. Часть результатов Велиева изложена (и частично улучшена) в книге автора [8]: некоторые его результаты усилены в статье Тополянского и Запрудского [1]. Отметим еще, что в работах А.З. Наматова [1, 2] исследована устоячивость метода Фаздо - Галеркина для задач, в которых квазилинейное параболическое уравнение решается при краевом условии (вообще говоря, нелинейном), содержащим производную по времени от искомого функции. Атакишева [1] и Запрудский [1] изучали устоячивость метода Бубнова - Галеркина в банаховых пространствах.

§ 6. Погрешность изображения рекуррентного вычислительного процесса.

1°. Рассмотрим бесконечный рекуррентный вычислительный процесс

$$n = 0, 1, 2, \dots; \sum_{k=0}^n A_{nk} x^{(k)} = f^{(n)}. \quad (Ш.6.1)$$

Примем, что $x^{(n)} \in X_n$, $f^{(n)} \in Y_n$, где X_n и Y_n - банаховы пространства, и что A_{nk} - линейные операторы, действующие из X_k в Y_n . Будем считать, что каждый из операторов A_{nk} ограниченно обратим, а произведения $A_{nn}^{-1} A_{nk}$, $0 \leq k \leq n-1$, ограничены. При этих предположениях существует последовательность элементов $x^{(n)} \in X_n$, удовлетворяющих бесконечной системе (Ш.6.1).

Пусть каждый из операторов A_{nk} вычислен с погрешностью Γ_{nk} , а каждый из элементов $f^{(n)}$ - с погрешностью $\delta^{(n)}$, так что вместо системы (Ш.6.1) на самом деле решается система

$$n \geq 0; \sum_{k=0}^n (A_{nk} + \Gamma_{nk}) x^{(k)} = f^{(n)} + \delta^{(n)}. \quad (Ш.6.2)$$

Допустим, что погрешности Γ_{nk} можно сделать сколь угодно малыми в следующем смысле: назовем бы нам (любо число $\varepsilon > 0$, можно добиться того, чтобы

$$\forall n \geq 0, \sum_{k=0}^n \|A_{nk}^{-1} \Gamma_{nk}\| \leq \varepsilon. \quad (Ш.6.3)$$

2°. Оценим погрешность искажения $\hat{\xi}_n = \|x^{(n)} - x^{(n)}\|$. Уравнение (Ш.6.1) и (Ш.6.2) преобразуем к виду

$$A_{nn} x^{(n)} = q^{(n)} := f^{(n)} - \sum_{k=0}^{n-1} A_{nk} x^{(k)}; \quad (\text{Ш.6.4})$$

$$(A_{nn} + \Gamma_{nn}) x^{(n)} = q^{(n)} + \varepsilon^{(n)};$$

$$x^{(n)} := \delta^{(n)} - \sum_{k=0}^{n-1} A_{nk} (x^{(k)} - x^{(k)}) - \sum_{k=0}^{n-1} \Gamma_{nk} x^{(k)}. \quad (\text{Ш.6.5})$$

Переход к этим новым уравнениям преобразует рекуррентный процесс в эквивалентный свободный, и для оценки погрешности искажения можно воспользоваться неравенством (Ш.1.7), которое в данном случае означает следующее: если $\|A_{nn}^{-1} \Gamma_{nn}\| \leq \beta$, $0 < \beta < 1$, то

$$\begin{aligned} \hat{\xi}_n &\leq \frac{1}{1-\beta} \left[\|A_{nn}^{-1} \Gamma_{nn} x_*^{(n)}\| + \|A_{nn}^{-1} \varepsilon^{(n)}\| \right] \leq \\ &\leq \frac{1}{1-\beta} \left[\|A_{nn}^{-1} \Gamma_{nn}\| \cdot \|x_*^{(n)}\| + \sum_{k=0}^{n-1} \|A_{nn}^{-1} \Gamma_{nk}\| \cdot \|x_*^{(k)}\| + \right. \\ &\quad \left. + \sum_{k=0}^{n-1} \left[\|A_{nn}^{-1} A_{nk}\| + \|A_{nn}^{-1} \Gamma_{nk}\| \right] \hat{\xi}_k + \|A_{nn}^{-1}\| \cdot \|\delta^{(n)}\| \right] = \\ &= \frac{1}{1-\beta} \left\{ \sum_{k=0}^n \|A_{nn}^{-1} \Gamma_{nk}\| \cdot \|x_*^{(k)}\| + \|A_{nn}^{-1}\| \cdot \|\delta^{(n)}\| + \sum_{k=0}^{n-1} \left[\|A_{nn}^{-1} A_{nk}\| + \|A_{nn}^{-1} \Gamma_{nk}\| \right] \hat{\xi}_k \right\} \end{aligned} \quad (\text{Ш.6.6})$$

Формула (Ш.6.6) позволяет последовательно вычислить оценки для $\hat{\xi}_0, \hat{\xi}_1, \hat{\xi}_2, \dots$, если известны оценки для $\|x_*^{(n)}\|$, $n = 0, 1, 2, \dots$.

§ 7. Устойчивость рекуррентного вычислительного процесса Γ^0 . Процесс (Ш.6.1) назовем устойчивым относительно искажения в последовательности пар пространств (X_n, Y_n) , если существуют такие постоянные $\rho_1, \rho_2, \rho_3 > 0$, что из неравенства

$$\forall n \geq 0, \sum_{k=0}^n \|A_{nn}^{-1} \Gamma_{nk}\| \leq \rho_3 \quad (\text{Ш.7.1})$$

вытекает следующая оценка погрешности искажения

$$\hat{\xi}_n \leq \rho_1 \max_{j < n} \sum_{k=0}^j \|A_{jj}^{-1} \Gamma_{jk}\| + \rho_2 \max_{j < n} \|A_{jj}^{-1} \delta^{(j)}\|. \quad (\text{Ш.7.2})$$

Это определение соответствует первому из определений устойчивости свободного вычислительного процесса.

Замечание. В предшествующем параграфе было отмечено, что любой рекуррентный вычислительный процесс легко преобразуется в свободный. Можно было бы определить устойчивость рекуррентного процесса как устойчивость соответствующего свободного процесса, и тогда были бы применимы утверждения § 2 настоящей главы. Однако, условия этих утверждений сказываются в этом случае громоздкими и трудно проверяемыми, почему мы и предпочитаем дать для рекуррентного процесса независимое определение устойчивости.

Теорема Ш.7.1. Пусть $h^{(n)} \in Y_n$, $n = 0, 1, 2, \dots$, и пусть $t^{(n)}$, $n = 0, 1, 2, \dots$, есть последовательность решений бесконечной рекуррентной системы

$$\forall n \geq 0, \sum_{k=0}^n A_{nk} t^{(k)} = h^{(n)}. \quad (\text{Ш.7.3})$$

Если выполнены неравенства

$$\|t^{(n)}\| \leq C_0 \max_{j < n} \|A_{jj}^{-1} h^{(j)}\|, \quad C_0 = \text{const} \quad (\text{Ш.7.4})$$

$$\|x_*^{(n)}\| \leq C_1 = \text{const}, \quad (\text{Ш.7.5})$$

где $x_*^{(n)}$, $n = 0, 1, 2, \dots$ - последовательность решений системы (Ш.6.1), то вычислительный процесс (Ш.6.1) устойчив.

Уравнение (Ш.6.2) преобразуем к виду

$$\sum_{k=0}^n A_{nk} (x^{(k)} - x_*^{(k)}) = \delta^{(n)} - \sum_{k=0}^n \Gamma_{nk} (x^{(k)} - x_*^{(k)}) - \sum_{k=0}^n \Gamma_{nk} x_*^{(k)}.$$

По неравенствам (Ш.7.4) и (Ш.7.5)

$$\hat{\xi}_n \leq C_0 \left\{ \max_{j < n} \|A_{jj}^{-1} \delta^{(j)}\| + \max_{j < n} \sum_{k=0}^j \|A_{jj}^{-1} \Gamma_{jk}\| \times \right. \\ \left. \times \max_{j < n} \hat{\xi}_j + C_1 \max_{j < n} \sum_{k=0}^j \|A_{jj}^{-1} \Gamma_{jk}\| \right\}. \quad (\text{Ш.7.6})$$

Обозначим

$$\sigma_n = \max_{j < n} \|A_{jj}^{-1} \delta^{(j)}\| + C_1 \max_{j < n} \sum_{k=0}^j \|A_{jj}^{-1} \Gamma_{jk}\|. \quad (\text{Ш.7.7})$$

Очевидно, σ_n не убывает с возрастанием n . Возьмем ρ_3 отоль малым, чтобы $\beta = C_0 \rho_3 < 1$. Тогда

$$\hat{\xi}_n \leq \beta \max_{j < n} \hat{\xi}_j + \sigma_n. \quad (\text{Ш.7.8})$$

Докажем, что

$$\forall n > 0, \hat{\xi}_n \leq \frac{\sigma_n}{1-\beta}. \quad (\text{Ш.7.9})$$

Неравенство (Ш.7.9) верно при $n=0$ - в этом случае (Ш.7.8) принимает вид $\hat{\xi}_0 \leq \beta \hat{\xi}_0 + \sigma_0$. Пусть неравенство (Ш.7.9) верно для $n < N-1$; докажем, что оно тогда верно и для $n=N$. Может случиться, что

$$\max_{j < N} \hat{\xi}_j = \hat{\xi}_N, \quad (\text{Ш.7.10})$$

тогда по (Ш.7.8), $\hat{\xi}_N \leq \beta \hat{\xi}_N + \sigma_N$, и неравенство (Ш.7.9) доказано. Если же неравенство (Ш.7.10) неверно, то $\max_{j < N} \hat{\xi}_j = \hat{\xi}_K$, где K - одно из чисел $0, 1, \dots, N-1$, и тогда

$$\hat{\xi}_N < \hat{\xi}_K \leq \frac{\sigma_K}{1-\beta} \leq \frac{\sigma_N}{1-\beta},$$

что и требовалось доказать. Итого, что теперь неравенство (Ш.7.2) справедливо при значениях постоянных

$$\rho_1 = \frac{C_1}{1-\beta}, \quad \rho_2 = \frac{1}{1-\beta}, \quad \rho_3 = \frac{\beta}{C_0}; \quad \forall \beta \in (0, 1). \quad (\text{Ш.7.11})$$

2°. Нетрудно указать простое условие, при котором неравенство (Ш.7.4) справедливо: достаточно, чтобы

$$\sum_{k=0}^{n-1} \|A_{nn}^{-1} A_{nk}\| \leq q = \text{const} < 1. \quad (\text{Ш.7.12})$$

Действительно, из уравнения (Ш.7.3) следует, что

$$\|t^{(n)}\| \leq q \max_{k < n-1} \|t^{(k)}\| + \|A_{nn}^{-1} h^{(n)}\|. \quad (\text{Ш.7.13})$$

Пусть $\max_{k < n} \|t^{(k)}\| = \|t^{(l)}\|$, $l \leq n$. Если $l = n$, то из (Ш.7.13) находим $\|t^{(n)}\| \leq q \|t^{(n)}\| + \|A_{nn}^{-1} h^{(n)}\|$, или $\|t^{(n)}\| \leq$

$\leq (1-q)^{-1} \|A_{nn}^{-1} h^{(n)}\|$. Если же $l < n$, то, заменив в (Ш.7.13) n на l , получаем $\|t^{(l)}\| \leq q \|t^{(l)}\| + \|A_{ll}^{-1} h^{(l)}\|$ и, следова-

тально,

$$\|t^{(n)}\| < \|t^{(l)}\| \leq (1-q)^{-1} \|\bar{A}_{nn}^{-1} h^{(l)}\| \leq (1-q)^{-1} \max_{j \leq n} \|\bar{A}_{jj} h^{(j)}\|.$$

3°. Теорема Ш.7.2. Рекуррентный процесс (Ш.6.1) устойчив, если выполнены условия (Ш.7.5) и (Ш.7.12).

Доказательство следует из неравенства (Ш.6.6). Допуская, что

$$\sum_{k=0}^n \|\bar{A}_{nn}^{-1} \Gamma_{nk}\| \leq \beta, \quad 0 < \beta < 1,$$

получаем из (Ш.6.6)

$$\hat{\xi}_n \leq \frac{1}{1-\beta} \left[C_1 \sum_{k=0}^n \|\bar{A}_{nn}^{-1} \Gamma_{nk}\| + \|\bar{A}_{nn}^{-1} \delta^{(n)}\| + (q+\beta) \max_{k \leq n} \hat{\xi}_k \right].$$

Выберем β столь малым, чтобы $q_1 := (q+\beta)(1-\beta)^{-1} < 1$. Имеем

$$\hat{\xi}_n \leq q_1 \max_{k < n} \hat{\xi}_k + \frac{1}{1-\beta} \left[C_1 \sum_{k=0}^n \|\bar{A}_{nn}^{-1} \Gamma_{nk}\| + \|\bar{A}_{nn}^{-1} \delta^{(n)}\| \right]. \quad (\text{Ш.7.14})$$

Уравнения (Ш.7.14) и (Ш.7.8) с точностью до обозначений совпадают, и аналогично формуле (Ш.7.9) получаем

$$\hat{\xi}_n \leq \rho'_1 \sum_{k=0}^n \|\bar{A}_{nn}^{-1} \Gamma_{nk}\| + \rho'_2 \|\bar{A}_{nn}^{-1} \delta^{(n)}\|. \quad (\text{Ш.7.15})$$

Замечание 2. Из замечания 1 (см. п. 1°) ясно, что теорема Ш.7.2 является частным случаем теоремы Ш.7.1. Мы считаем целесообразным выделить этот частный случай, потому что оценки (Ш.7.15) несколько лучше по порядку, чем это требуется формулой (Ш.7.2); входящей в определение устойчивости.

§ 8. Точность искажения и устойчивость метода наискорейшего спуска.

1°. Результаты §§ 6,7 применим к методу наискорейшего спуска, предложенному в свое время Дж. Темплом [1] и, независимо от него, Канторовичем [2] (см. также Канторович и Акилов [1]), который установил также оценку скорости сходимости метода (см. ниже, формула (Ш.8.5)); в статье Темпля установлен только самый факт сходимости метода. Коротко напомним основные черты метода наискорейшего спуска.

Пусть \bar{A} - ограниченный самосопряженный положительно определенный оператор, действующий в некотором гильбертовом простран-

отве H , и пусть потребуе решить уравнение

$$Ax = f. \quad (Ш.8.1)$$

Зададим произвольным элементом $x^{(0)} \in H$, который будем рассматривать как нулевое приближение к искомому решению x_* уравнения (Ш.8.1). Пусть построено $(n-1)$ -е приближение $x^{(n-1)}$, тогда следующее приближение $x^{(n)}$ строится по формулам

$$y^{(n-1)} = Ax^{(n-1)} - f, \quad \sigma_{n-1} = \frac{\|y^{(n-1)}\|^2}{(Ay^{(n-1)}, y^{(n-1)})}, \quad x = x^{(n-1)} - \sigma_{n-1} y^{(n-1)}. \quad (Ш.8.2)$$

Как показано в статье Канторовича [2], $x^{(n)} \xrightarrow{n \rightarrow \infty} x_*$ с оценкой

$$\|x_* - x^{(n)}\| \leq \frac{\|y_1\|}{m} \left(\frac{M-m}{M+m} \right)^n, \quad (Ш.8.3)$$

где M и m суть соответственно верхняя и нижняя грани оператора A :

$$M = \sup_{\|f\|=1} (Af, f), \quad m = \inf_{\|f\|=1} (Af, f).$$

Таким образом, метод наискорейшего спуска связан с бесконечным рекуррентным процессом

$$x^{(n)} - B_n x^{(n-1)} = f^{(n)}, \quad (Ш.8.4)$$

где

$$B_n = I - \sigma_{n-1} A, \quad f^{(n)} = \sigma_{n-1} f. \quad (Ш.8.5)$$

2^o Исследуем погрешность искажения и устойчивость процесса (Ш.8.4). Очевидно, $M \leq \sigma_{n-1} \leq m$; верхняя и нижняя грани оператора B_n суть, следовательно, $1 - \sigma_{n-1} m$ и $1 - \sigma_{n-1} M$.

Далее,

$$1 - \sigma_{n-1} M > \frac{M-m}{m}, \quad 1 - \sigma_{n-1} m < \frac{M-m}{M}.$$

Отсюда

$$\|B_n\| \leq \max \left(\frac{M-m}{m}, \frac{M-m}{M} \right) = \frac{M-m}{m}. \quad (Ш.8.6)$$

В обозначениях § 6 $A_{nn} = I$, $A_{n,n-1} = -B_n$, $A_{nk} = 0$, $k \leq n-2$.

Обозначим через δ_n погрешность вычисления величины σ_{n-1} , тогда

$\Gamma_{n,n-1} = \delta_n A$ и $\Gamma_{nk} = 0$ при $k \neq n-1$; кроме того $\delta_n^{(n)} = \delta_n f$.

Обратимся к формуле (Ш.8.6). Условие $\|A_{nn} \Gamma_{nn}\| \leq \beta$, $0 < \beta < 1$,

здесь выполняется тривиально, потому что $\Gamma_{mn} = 0$. Мож. считать $\beta = 0$, и формула (Ш.6.5) принимает вид

$$\hat{\xi}_n \leq \left(\frac{M-m}{m} + |\vartheta_n| M \right) \hat{\xi}_{n-1} + |\vartheta_n| (C_0 M + \|f\|). \quad (\text{Ш.8.7})$$

Здесь C_0 - какая-нибудь верхняя граница значений $\|x_n^{(0)}\|$; такая конечная граница существует, потому что процесс наискорейшего спуска охотится. Эту границу нетрудно вычислить, исходя из неравенства (Ш.8.3).

Проанализируем формулу (Ш.8.7). Пусть погрешность ϑ_n достаточно мала, $|\vartheta_n| \leq \varepsilon$. Если при этом $M < 2m$, то

$$\frac{M-m}{m} + |\vartheta_n| M \leq \frac{M-m}{m} + \varepsilon M =: q$$

и при ε достаточно малом будет $q < 1$. Теперь

$$\hat{\xi}_n \leq q \hat{\xi}_{n-1} + \varepsilon a, \quad a = C_0 M + \|f\|. \quad (\text{Ш.8.8})$$

Заметим еще, что

$$\hat{\xi}_0 = 0 \quad (\text{Ш.8.9})$$

потому что элемент $x^{(0)}$ не вычисляется, а задается, и соответственно считать, что $x^{(0)} = x^{(0)}$. Очевидно, что $\hat{\xi}_n \leq \tau_n$, где τ_n - решение разностного уравнения

$$\tau_n = q \tau_{n-1} + \varepsilon a, \quad \tau_0 = 0. \quad (\text{Ш.8.10})$$

Решение этого уравнения есть

$$\tau_n = \frac{\varepsilon a}{1-q} (1-q^n),$$

и оценка погрешности искажения имеет вид

$$\hat{\xi}_n \leq \frac{\varepsilon a}{1-q} (1-q^n) < \frac{\varepsilon a}{1-q}. \quad (\text{Ш.8.11})$$

Таким образом, если: 1) $M < 2m$ и 2) ε - верхняя граница погрешности величины σ_{n-1} - достаточно мала, то погрешность искажения метода наискорейшего спуска есть величина порядка $O(\varepsilon)$, и эта оценка - равномерная относительно n .

Если хотя бы одно из условий 1), 2) нарушено, то может случиться, что погрешность искажения будет неограниченно возрастать вместе с n .

5°. Близкие результаты дает и анализ устойчивости. Условие (Ш.7) очевидно выполнено, а условие (Ш.7.12) в данном случае означает, что $\|B_n\| \leq q < 1$; это условие будет выполнено, если $M < 2m$ - тогда можно положить $q = (M-m)/m$. По теореме

(Ш.7.2) процессом наискорейшего спуска устойчив, если $M < 2m$. Тот же результат дает и теорема Ш.7.1.

4°. Если $M \geq 2m$, то для оценки погрешности можно воспользоваться непосредственно формулой (ш.8.7). Допустим, что величина $\hat{\sigma}_n$ ограничена, тогда $(M-m)/m + \hat{\sigma}_n M \leq \tau = \text{const}$ и

$$\hat{\xi}_n \leq \tau \hat{\xi}_{n-1} + |\hat{\sigma}_n| a, \quad \hat{\xi}_0 = 0. \quad (\text{Ш.8.12})$$

Решение этого неравенства есть

$$\hat{\xi}_n \leq a \sum_{k=1}^n |\hat{\sigma}_k| \tau^{n-k} \quad (\text{Ш.8.13})$$

Зафиксируем натуральное число N . Если желательно, чтобы было $\hat{\xi}_k \leq \bar{\xi}$, где $\bar{\xi}$ — заданное число, то достаточно выбрать $\hat{\sigma}_{k-1}$ для $1 \leq k \leq N$ с такой точностью, чтобы

$$|\hat{\sigma}_k| \leq \frac{\bar{\xi}}{N \tau^{N-k}}. \quad (\text{Ш.8.14})$$

§ 9. Об устойчивости метода коллокации.

1°. Пусть неискаженная алгебраическая система метода коллокации имеет вид

$$M_n a^{(n)} = f^{(n)}, \quad (\text{Ш.9.1})$$

где M_n — матрица чисел $\psi_{kl}(t_j^{(n)}) = (A \psi_{kl})(t_j^{(n)})$, $a^{(n)} = (a_1^{(n)}, a_2^{(n)}, \dots, a_r^{(n)})$, $f^{(n)} = (f(t_1^{(n)}), f(t_2^{(n)}), \dots, f(t_N^{(n)}))$.

Пусть \mathcal{K} — компакт в m -мерном евклидовом пространстве R_m . Точки $t_j^{(n)}$ — узлы коллокации — выберем в вершинах некоторой параллелепипедальной сетки, лежащих в \mathcal{K} . Ребра параллелепипедов сетки параллельны k -й координатной оси, пусть имеют длину $h_k^{(n)}$, $k = 1, 2, \dots, m$. Допустим, что $C_1 h_n \leq h_k^{(n)} \leq C_2 h_n$, где C_1 и C_2 — постоянные и $h_n \rightarrow 0$.

Будем рассматривать $a^{(n)}$ и элемент пространства R_r , а $f^{(n)}$ — как элемент пространства Y_N , также N -мерного, в котором норма определяется так:

$$\|f^{(n)}\|_{Y_N} = h_n^{-1/2} \|f^{(n)}\|_{R_N}, \quad (\text{Ш.9.2})$$

β — вещественная постоянная. Оператор, порожденный матрицей M_n и действующий из R_r в Y_N , обозначим через A_n ; систему (Ш.9.1) можно теперь записать в такой операторной форме:

$$A_n \alpha^{(n)} = f^{(n)}; \quad (Ш.9.3)$$

соответствующую искаженную систему запишем в обычной форме

$$(A_n + \Gamma_n) \beta^{(n)} = f^{(n)} + \delta^{(n)}. \quad (Ш.9.4)$$

Оператор, порожденный матрицей M_n и действующий в R_N , обозначим тем же символом M_n .

2^o. Пусть $\delta_n^{(k)}$, $k = 1, 2, \dots, N$ - сингулярные числа матрицы M_n , т.е., собственные числа неотрицательной матрицы $M_n^* M_n$. Расположим их в порядке возрастания: $\delta_n^{(1)} \leq \delta_n^{(2)} \leq \dots \leq \delta_n^{(N)}$.

Теорема Ш.9.1. Если

$$\delta_n^{(1)} \geq C_0 h_n^{-1}, \quad C_0 = \text{const} > 0 \quad (Ш.9.5)$$

и $\|a^{(n)}\| \leq C_1 = \text{const}$, то вычислительный процесс метода коллокации (Ш.9.3) устойчив в смысле второго определения (§ 2, гл. IV) в последовательности пар пространств (R_N, Y_{N3}) . Если же

$$\delta_n^{(1)} = r(h_n) h_n^{-1}, \quad r(h_n) \xrightarrow{h_n \rightarrow \infty} 0, \quad (Ш.9.6)$$

то в той же последовательности процесс (Ш.9.3) неустойчив.

По следствию Ш.2.2, для устойчивости процесса (Ш.9.3) в упомянутом смысле необходимо и достаточно, чтобы: 1) $\|A_n^{-1}\| \leq C_1$, 2) $\|A_n^{-1}\| \cdot \|a^{(n)}\|_{R_N} \leq C_2$, где C_1 и C_2 - постоянные. Оценим величину

$$\begin{aligned} \|A_n^{-1}\| &= \sup_{g^{(n)} \in Y_{N3}} \frac{\|A_n^{-1} g^{(n)}\|_{R_N}}{\|g^{(n)}\|_{Y_{N3}}} = \\ &= h_n^{-1/2} \sup_{g^{(n)} \in R_N} \frac{\|M_n^{-1} g^{(n)}\|_{R_N}}{\|g^{(n)}\|_{R_N}} = h_n^{-1/2} \|M_n^{-1}\|. \end{aligned}$$

Из соотношения $M_n^{-1} (M_n^*)^{-1} = (M_n^* M_n)^{-1}$ следует, что сингулярные числа матрицы M_n^{-1} суть $1/\delta_n^{(1)}, 1/\delta_n^{(2)}, \dots, 1/\delta_n^{(N)}$.

Отсюда $\|M_n^{-1}\| = \|(M_n^* M_n)^{-1}\|^{1/2} = 1/\sqrt{\delta_n^{(1)}}$ и окончательно

$$\|A_n^{-1}\| = \frac{1}{\sqrt{h_n \delta_n^{(1)}}}. \quad (Ш.9.7)$$

Если верно неравенство (Ш.9.6), то $\|A_n^{-1}\| \geq [r(h_n)]^{1/2} h_n \xrightarrow{h_n \rightarrow \infty} \infty$,

и процесс (Ш.9.3) неустойчив.

Пусть справедливо неравенство (Ш.9.5). Тогда $\|A_n^{-1}\| \leq C_0^{-1/2}$, и условие 1) выполнено. Условие 2) выполнено по предположению, и процесс коллокации устойчив.

Отметим частный случай. Прежде всего заметим следующее: так как $f \in C(\mathcal{X})$, то тем более $f \in L_2(\mathcal{X})$, а тогда

$$h_n^m \sum_{j=1}^N f_j^2 \leq C = \text{const}$$

или, что то же, $\|f^{(n)}\|_{y_{n,m}} \leq C$. Пусть теперь неравенство (Ш.9.5) выполнено при $\delta = m$, тогда норма $\|A_n^{-1}\|$ ограничена постоянной $C^{-1/2}$. Как только что было отмечено, нормы $\|f^{(n)}\|_{y_{n,m}}$ также ограничены, а тогда ограничены и нормы $\|a^{(n)}\|_{R_N} \leq \|A_n^{-1}\| \cdot \|f^{(n)}\|_{y_{n,m}}$.

По следствию Ш.2.2, процесс (Ш.9.3) устойчив. Таким образом, при $\delta = m$ условие (Ш.9.5) необходимо и достаточно для устойчивости процесса коллокации (Ш.9.3).

3°. Выше (см. § 3, гл.Ш) было дано определение сильной минимальности последовательности систем

$$(y_{n1}, y_{n2}, \dots, y_{nN}); \quad N = N(n), \quad n = 1, 2, \dots \quad (\text{Ш.9.8})$$

Введем еще одно определение: систему последовательностей (Ш.9.8) назовем почти ортонормированной, если собственные числа матриц A_n этих систем положительно ограничены сверху и снизу. Очевидно, почти ортонормированная последовательность также и сильно минимальна. Это определение обобщает известное определение Талдыкина [1]. Системы, которые мы называем почти ортонормированными, Талдыкин называет регулярными.

Теорема Ш.9.2. Пусть X - бесконечномерное гильбертово пространство, и последовательность (Ш.9.8) почти ортонормирована в X , а процесс (Ш.9.3) устойчив в смысле второго определения (§ 2, гл.Ш). Тогда процесс вычисления приближенного решения по методу коллокации устойчив в смысле того же определения в последовательности пар пространств (X_n, Y_n) . Если последовательность (Ш.9.8) сильно минимальна в X и процесс (Ш.9.3) неустойчив, то процесс вычисления приближенного решения также неустойчив.

Если последовательность (Ш.9.8) почти ортонормирована в X

то

$$\|x_n^{(n)} - \tilde{x}_n^{(n)}\|_X^2 = (M_n(a^{(n)} - b^{(n)}), a^{(n)} - b^{(n)}) \leq \Lambda_0 \|a^{(n)} - b^{(n)}\|_{R_N}^2;$$

здесь M_n - матрица Грама элементов (Ш.9.8) при фиксированном n , Λ_0 - верхняя граница собственных чисел этих матриц при всех n , $\alpha_x^{(n)}$ и $x_x^{(n)}$ - соответственно непокаженное и искаженное приближенные решения, с которыми векторы $a^{(n)}$ и $b^{(n)}$ связаны соотношениями $\alpha_x^{(n)} = \sum_{k=1}^n a_k^{(n)} y_{nk}$, $x_x^{(n)} = \sum_{k=1}^n b_k^{(n)} y_{nk}$. Так как процесс вычисления коэффициентов устойчив, то существуют такие постоянные $\rho, q, \nu > 0$, что при $\| \Gamma_n \| \leq \nu$ будет

$$\| x_x^{(n)} - x_x^{(n)} \| \leq \sqrt{\Lambda_0} (\rho \| \Gamma_n \|_{R_x} \rightarrow y_{R_x} + q \| \delta^{(n)} \|_{y_{R_x}});$$

последнее неравенство означает, что процесс вычисления приближенного решения устойчив.

Если последовательность (Ш.9.8) сильно минимальна в X , то, обозначая положительную нижнюю границу собственных чисел матрицы M_n через λ_0 , имеем

$$\| x_x^{(n)} - x_x^{(n)} \| \geq \sqrt{\lambda_0} \| a^{(n)} - b^{(n)} \|_{R_x}.$$

По предположению, процесс (Ш.9.8) неустойчив, поэтому при любых фиксированных независимо от n , нормах $\| \Gamma_n \|$ и $\| \delta^{(n)} \|$ нормы $\| a^{(n)} - b^{(n)} \|$ неограничены, а тогда неограничены и нормы $\| x_x^{(n)} - x_x^{(n)} \|$, и процесс вычисления приближенных решений неустойчив.

4°. Некоторые оценки наименьшего сингулярного числа $\delta_n^{(1)}$ можно получить, исследуя определитель $D_n = \det M_n$. Ограничимся случаем $m=1$, $\mathcal{K}=[0, 1]$, и пусть узлы $t_j^{(n)} = j/n$, $j=1, 2, \dots, n$. Допустим, что функции $\Psi_{nk}(t) = (A y_{nk})(t)$ удовлетворяют условию Липшица с фиксированным показателем β и с постоянной, которая оценивается некоторой степенью K :

$$| \Psi_{nk}(t') - \Psi_{nk}(t'') | \leq C K^\beta |t' - t''|^\beta; \quad \beta = \text{const}; \quad (\text{Ш.9.9})$$

где положим еще, что

$$| \Psi_{nk}(t) | \leq C K^d, \quad d = \text{const} \geq 0 \quad (\text{Ш.9.10})$$

Имеет (для упрощения обозначений пишем Ψ_k вместо Ψ_{nk} в t_j вместо $t_j^{(n)}$)

$$D_n = \begin{vmatrix} \psi_1(t_1) & \psi_2(t_1) & \dots & \psi_n(t_1) \\ \psi_1(t_2) & \psi_2(t_2) & \dots & \psi_n(t_2) \\ \dots & \dots & \dots & \dots \\ \psi_1(t_n) & \psi_2(t_n) & \dots & \psi_n(t_n) \end{vmatrix}. \quad (\text{Ш.9.1})$$

Из первой строки определителя D_n вычтем вторую, из второй - третью и т.д.; из $(n-1)$ -й строки вычтем n -ю. Элементы первых $n-1$ строк имеют вид $\psi_k(t_j) - \psi_k(t_{j+1})$ и оцениваются, по равенству (Ш.9.9), следующим образом:

$$|\psi_k(t_j) - \psi_k(t_{j+1})| \leq C n^{-r} \kappa^\beta;$$

элементы l следней строки имеют оценку $C \kappa^\alpha$. По теореме Адамса (см., например, Смирнов [1], стр. 64) получаем

$$|D_n| \leq \frac{C^n}{n^{n(n-1)}} \left[\sum_{k=1}^n \kappa \right]^{1/2} \left[\sum_{k=1}^n \kappa^{2\beta} \right]^{n-1/2}.$$

Так как $\alpha \geq 0$, то $\kappa^{2\alpha} \leq \int_{\kappa}^{\kappa+1} u^{2\alpha} du$ и

$$\sum_{k=1}^n \kappa^{2\alpha} < \int_1^{\kappa+1} u^{2\alpha} du < (2\alpha+1)^{-1} (n+1)^{2\alpha+1};$$

аналогично

$$\sum_{k=1}^n \kappa^{2\beta} < (2\beta+1)^{-1} (n+1)^{2\beta+1}.$$

Теперь

$$|D_n| \leq \frac{C_1^n}{n^{n(n-1)}} (n+1)^{\frac{2\alpha+1}{2} + \frac{(2\beta+1)(n-1)}{2}}; \quad C_1 = \text{const},$$

откуда легко получается и несколько более простая оценка

$$|D_n| \leq C_2^n; \quad C_2 = \text{const}.$$

Далее, $D_n^2 = \text{Det}(M_n^* M_n) = \dots \geq [S_n^{(n)}]^n;$ (Ш.9.2)

откуда

$$C_2 \leq D_n \leq C_3 n^{\frac{1}{2}(2\alpha+2r)(n-1) + \frac{2\alpha+1}{2}} \leq C_3 n^{(\beta-r)+1}; \quad C_3 = \text{const}. \quad (\text{Ш.9.3})$$

Теперь из теоремы (Ш.9.1) вытекает, что если $m = j = 1$ и $\beta > \gamma$, то процесс (Ш.9.3) неустойчив. В частности, если $\gamma = 1$, то процесс (Ш.9.3) неустойчив при $\beta > 1$.

5°. Теорема Адамара часто дает грубую оценку определителя. Рассмотрим пример, в котором эту оценку можно улучшить.

Пусть A - тождественный оператор. Задача состоит тогда в интерполяции функции $f(t)$ линейной комбинацией вида

$$\sum_{k=1}^n a_k^{(n)} y_{nk}(t)$$

на заданном компакте \mathcal{K} . Пусть $m = 1$, $\mathcal{K} = [0, 1]$ и $y_{nk}(t) = y_k(t) = t^{k-1}$, $k = 1, 2, \dots$. Положим еще $t_j = t_j^{(0)} = j/n$, $j = 0, 1, 2, \dots, n-1$. Определитель коллокационной системы

$$D_n = \begin{vmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & t_1 & t_2 & \dots & t_{n-1} \\ 1 & t_1^2 & t_2^2 & \dots & t_{n-1}^2 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & t_1^{n-1} & t_2^{n-1} & \dots & t_{n-1}^{n-1} \end{vmatrix} = t_1 t_2 \dots t_{n-1} \begin{vmatrix} 1 & 1 & \dots & 1 \\ t_1 & t_2 & \dots & t_{n-1} \\ \dots & \dots & \dots & \dots \\ t_1^{n-2} & t_2^{n-2} & \dots & t_{n-1}^{n-2} \end{vmatrix}$$

равен определителю Вандермонда чисел t_1, t_2, \dots, t_{n-1} , умноженному на $t_1 t_2 \dots t_{n-1} = \frac{(n-1)!}{n^{n-1}}$. Таким образом,

$$D_n = \frac{(n-1)!}{n^{n-1}} \prod_{1 \leq k < j \leq n-1} (t_j - t_k) < 1.$$

В таком случае $\rho_n^{(1)} < 1 = h \cdot h^{-1}$; $h = 1/n$, $r(h) = h \rightarrow 0$

и коллокационный процесс неустойчив в последовательности нар пространств (R_n, Y_{n2}) .

6°. В заключении параграфа коротко изложим содержание работы Вайншико [8], в которой, по-видимому, впервые рассмотрен вопрос об устойчивости метода коллокации. В этой работе рассматривается обыкновенное дифференциальное уравнение

$$x^{(m)}(t) - \sum_{k=1}^{m-1} c_k(t) x^{(k)}(t) = f(t) \tag{Ш.9.14}$$

при краевых условиях

$$\sum_{k=0}^{m-1} [d_{jk} x^{(k)}(a) + \beta_{jk} x^{(k)}(b)] = 0; \tag{Ш.9.15}$$

$$j = 0, 1, \dots, m-1; d_{jk}, \beta_{jk} = \text{const.}$$

Предположим, что соответствующая однородная задача имеет только тривиальное решение. Тогда существует последовательность полиномов степени $m+k$

$$p_k(t) = \sum_{\ell=0}^{m+k} c_{k\ell} t^\ell; \quad k = (1, 2, \dots,$$

удовлетворяющая краевым условиям (Ш.9.15). Эту последовательность примем за координатную систему метода коллокации. За узлы коллокации при данном n примем корни $t_j^{(n)}$ полинома $\omega_n(t)$, где $\{\omega_n(t)\}$ - система полиномов, ортонормированная и полная в метрике $L_2(a, b; \rho(t))$, причем вес $\rho(t)$ есть функция, неотрицательная и суммируемая на (a, b) , такая, что

$$\int_a^b [\rho(t)]^{-1} dt < \infty.$$

Таким образом, приближенное решение задачи ищется как функция вида

$$x_n(t) = \sum_{k=0}^{n-1} a_k^{(n)} p_k(t); \quad (\text{Ш.9.16})$$

коэффициенты $a_k^{(n)}$ определяются из системы уравнений

$$\sum_{k=0}^{n-1} [P_k^{(m)}(t_j^{(n)}) - \sum_{\ell=0}^{m-1} c_{j\ell}^{(n)} p_k(t_j^{(n)})] a_k^{(n)} = f(t_j^{(n)}); \quad j=1, 2, \dots, n. \quad (\text{Ш.9.17})$$

В предположении, что коэффициенты $c_{k\ell}(t)$ и свободный член $f(t)$ непрерывны, доказываем, что система (Ш.9.17) однозначно разрешима при достаточно больших n ; приближенные решения $x_n(t)$ сходятся к точному решению $x_*(t)$ в метрике $L_2(a, b; \rho(t))$. Потребность аппроксимации имеет оценку

$$\|x_n^{(n)} - x_n^{(m)}\|_{L_2(a, b; \rho)} \leq \delta_n(x_n^{(m)}), \quad (\text{Ш.9.18})$$

где $\delta_n(u)$ есть наилучшее приближение функции $u(t)$ полиномами степени не выше n в метрике $L_2(a, b; \rho)$.

Обозначим через H_n подпространство пространства $L_2(a, b; \rho)$, образованное полиномами вида (Ш.9.16) с произвольными коэффициентами $a_k^{(n)}$. Доказываем, что при определенном выборе приближенных решений $x_n(t)$ устойчива в нормальном пространстве (H_n, R_n) тогда и только тогда, когда последовательность $\{p_k^{(m)}(t)\}$ сильно ϵ -стационарна в $L_2(a, b; \rho)$.

ГЛАВА IV

Погрешности алгоритма и округления

§ I. Число обусловленности.

1°. Пусть A — ограниченный и ограниченно обратимый оператор, действующий из банахова пространства X в такое же пространство Y . Числом обусловленности оператора A называется произведение

$$P_A = \|A\| \cdot \|A^{-1}\|. \quad (IV.1.1)$$

Очевидно, $P_A \geq 1$.

Если $X = Y = R_n$, то A есть числовая матрица порядка $n \times n$. Если эта матрица положительно определенная, то $\|A\| = \lambda_n^{(n)}$, $\|A^{-1}\| = 1/\lambda_n^{(1)}$, где $\lambda_n^{(1)}$ и $\lambda_n^{(n)}$ суть собственные числа матрицы A , соответственно наименьшее и наибольшее; в этом случае число обусловленности равно

$$P_A = \lambda_n^{(n)} / \lambda_n^{(1)}. \quad (IV.1.2)$$

Если A — произвольная эрмитова матрица в R_n , то матрица A^*A положительно определенная. Пусть $\delta_n^{(n)}$ и $\delta_n^{(1)}$ — наибольшее и наименьшее собственные числа последней матрицы, иначе говоря, $\delta_n^{(n)}$ и $\delta_n^{(1)}$ суть наибольшее и наименьшее сингулярные числа матрицы A . Тогда

$$P_A = \sqrt{\delta_n^{(n)} / \delta_n^{(1)}}. \quad (IV.1.3)$$

2°. Число обусловленности оператора имеет большое значение в некоторых случаях, касающихся оценки погрешности; один такой случай отмечен выше, в § 4 гл. I. Укажем еще один важный случай. Пусть A — ограниченный и ограниченно обратимый оператор, действующий в некотором сепарабельном гильбертовом пространстве H . Уравнение

$$Ax = f \quad (IV.1.4)$$

будем приближенно решать по методу наименьших квадратов: каждому натуральному n сопоставим другое натуральное N и выделим из H подпространство N измерений $H^{(n)}$; будем, далее, искать такой элемент $x_*^{(n)} \in H^{(n)}$, чтобы $\|Ax_*^{(n)} - f\| = \min$. Существование такого элемента очевидно. Оценим (см. статью автора и Р.А. Радевой [1]) погрешность аппроксимации $\rho_n = \|x_* - x_*^{(n)}\|$. Имеем

$$x_* - x_*^{(n)} = A^{-1} A(x_* - x_*^{(n)}) = A^{-1}(f - Ax_*)$$

и, следовательно, $\rho_n \leq \|A^{-1}\| \cdot \|f - Ax_*^{(n)}\|$. Если $x^{(n)}$ - произвольный элемент из $H^{(n)}$, то

$$\|f - Ax_*^{(n)}\| \leq \|f - Ax^{(n)}\| = \|A(x_* - x^{(n)})\| \leq \|A\| \cdot \|x_* - x^{(n)}\|; \rho_n \leq P_A \|x_* - x^{(n)}\|.$$

Взяв справа нижнюю грань по всевозможным $x^{(n)} \in H^{(n)}$ и обозначив

$$\forall x \in H, \varepsilon_n(x) = \inf_{x^{(n)} \in H^{(n)}} \|x - x^{(n)}\|,$$

получим интересующую нас оценку:

$$\rho_n = \|x_* - x_*^{(n)}\| \leq P_A \varepsilon_n(x_*). \quad (IV.1.5)$$

§ 2. Погрешность алгоритма в итерационном процессе

1°. Как хорошо известно, если A - самосопряженный оператор в гильбертовом пространстве H и его нижняя и верхняя грани удовлетворяют неравенствам $0 < m \leq M < \infty$, то задачу (IV.1.4) можно легко преобразовать так, чтобы она решалась итерациями. Достаточно, именно, заменить уравнение (IV.1.4) таким:

$$(I - B)x = \frac{2}{M+m}f; \quad B = I - \frac{2}{M+m}A. \quad (IV.2.1)$$

Нижняя и верхняя грани оператора B равны соответственно $-(M-m)/(M+m)$, поэтому

$$\|B\| = \frac{M-m}{M+m} = \frac{P_A - 1}{P_A + 1} < 1 \quad (IV.2.2)$$

и поэтому итерации для уравнения (IV.2.2) сходятся как прогрессия с знаменателем $(P_A - 1)/(P_A + 1)$. Ясно, что сходимость тем быстрее, чем меньше число обусловленности; если P_A достаточно велико, то сходимость итераций будет сколь угодно медленной.

2°. Выказанные в п.1° соображения применим к классическому методу Рунца. В separable гильбертовом пространстве H_0 рассмотрим задачу о минимуме функционала $|x|^2 - 2lx$, где $| \cdot |$ - норма в H_0 , а l - линейный функционал, определенный на всем пространстве H_0 и ограниченный в нем. Выберем в H_0 координатную систему $\{y_n\}$ и обозначим через $H_0^{(n)}$ подпространство пространства H_0 с базисом (y_1, y_2, \dots, y_n) . Приближенное решение $x_*^{(n)}$ нашей вариационной задачи построим как элемент подпространства $H_0^{(n)}$, реализующий минимум функционала $|x|^2 - 2lx$ на этом подпространстве. Очевидно,

$$x_*^{(n)} = \sum_{k=1}^n a_k^{(n)} y_k;$$

коэффициенты $a_k^{(n)}$ определяются из алгебраической системы

$$\sum_{k=1}^n [y_k, y_j] a_k^{(n)} = l y_j; \quad j = 1, 2, \dots, n,$$

в которой квадратные скобки означают скалярное произведение в H_0 . Последнюю систему запишем в виде векторного уравнения

$$M_n a^{(n)} = f^{(n)}; \tag{IY.2.3}$$

смысл обозначений очевиден. Обозначим через $\lambda_n^{(j)}$ и $\lambda_n^{(n)}$ наименьшее и наибольшее собственные числа матрицы Рунца M_n и положим

$$B_n = I_n - 2M_n / (\lambda_n^{(j)} + \lambda_n^{(n)}).$$

Уравнение (IY.2.3) заменим равносильным

$$a^{(n)} - B_n a^{(n)} = \frac{2f^{(n)}}{\lambda_n^{(j)} + \lambda_n^{(n)}}, \tag{IY.2.4}$$

которое разрешимо итерациями, так как $\|B_n\| = (P_n - 1)/(P_n + 1)$

$P_n = P_{M_n}$. Если начальное приближение равно нулю, и если ограничиться N итерациями, то соответствующая погрешность алгоритма для уравнения (IY.2.4) имеет оценку

$$\frac{2 \|f^{(n)}\| \cdot \|B_n\|^{N+1}}{\lambda_n^{(j)} + \lambda_n^{(n)}} = \frac{\|f^{(n)}\| \cdot \|B_n\|^{N+1}}{\lambda_n^{(j)} + \lambda_n^{(n)}}. \tag{IY.2.5}$$

Если $\lambda_n^{(j)} \rightarrow 0$ или $\lambda_n^{(n)} \rightarrow \infty$, то $P_n \rightarrow \infty$, $\|B_n\| \rightarrow 1$ и оценка погрешности алгоритма ухудшается. Желательно поэтому использовать такую координатную систему, для которой $\lambda_0 \leq \lambda_n \leq \lambda_n^{(n)} \leq \Lambda_0$, где λ_0 и Λ_0 - положительные постоянные; иначе говоря, целесообразно в качестве координатных брать почти ортонормированную систему в H_0 . В этом случае можно упростить определение матрицы B_n , именно, можно положить $B_n = I_n - 2M_n / (\Lambda_0 + \lambda_0)$ и вместо (IY.2.4) составить уравнение

$$a^{(n)} - B_n a^{(n)} = \frac{2f^{(n)}}{\Lambda_0 + \lambda_0}, \tag{IY.2.5'}$$

Такое равносильное уравнение (IY.2.5'). При таком определении $\|B_n\|$ и значения грани матрицы B_n соответственно равны

$$1 - \frac{2\lambda_n^{(n)}}{\Lambda_0 + \lambda_0} \geq 1 - \frac{2\Lambda_0}{\Lambda_0 + \lambda_0} = - \frac{\Lambda_0 - \lambda_0}{\Lambda_0 + \lambda_0}$$

и

$$1 - \frac{2\lambda_n^{(n)}}{\lambda_0 + \lambda_0} < 1 - \frac{2\lambda_0}{\lambda_0 + \lambda_0} = \frac{\lambda_0 - \lambda_0}{\lambda_0 + \lambda_0}.$$

Отсюда $\|B_n\| \leq (\lambda_0 - \lambda_0) / (\lambda_0 + \lambda_0)$; правая часть этого неравенства постоянна и меньше единицы. Для погрешности алгоритма после N итераций получаем оценку

$$\frac{\|f^{(n)}\|}{\lambda_0} \left(\frac{\lambda_0 - \lambda_0}{\lambda_0 + \lambda_0} \right)^{N+1}. \quad (IV.2.6)$$

Докажем, что $\|f^{(n)}\| \leq C$, где C не зависит от n . Почти нормированная система сильно минимальна; в книге автора [8] доказано (см. теорему 7.2), что если координатная система сильно минимальна, то последовательность векторов $a^{(n)}$, дополненная нулями, имеет предел в l_2 . Отсюда вытекает, что последовательность норм $\|a^{(n)}\|_{R_n}$ ограничена. Пусть $\|a^{(n)}\|_{R_n} \leq C_0 = \text{const}$.

Тогда из уравнения (IV.2.3) получаем

$$\|f^{(n)}\| \leq \|M_n\| \cdot \|a^{(n)}\| = \lambda_n^{(n)} \cdot \|a^{(n)}\| \leq \lambda_0 C_0.$$

Теперь из (IV.2.6) следует оценка погрешности алгоритма, не зависящая от n :

$$\frac{\lambda_0 C_0}{\lambda_0} \left(\frac{\lambda_0 - \lambda_0}{\lambda_0 + \lambda_0} \right)^{N+1}. \quad (IV.2.7)$$

Соображения и результаты настоящего параграфа легко переносятся на обобщенный метод Рунца.

§ 3. Погрешность округления в рекуррентном процессе

Γ^0 . Обратимся к системе (III.6.2); будем считать, что она может быть как конечной, так и бесконечной. Для дальнейшего vamos операторы $A_{nk} + \Gamma_{nk}$ и свободные члены $f + \delta^{(n)}$ известны. Для краткости обозначим

$$A_{nk} + \Gamma_{nk} = \hat{A}_{nk}, \quad f + \delta^{(n)} = \hat{f}_n, \quad (IV.3.1)$$

тогда система (III.6.2) запишется в виде

$$\forall n \geq 0, \quad \sum_{k=0}^r \hat{A}_{nk} z^{(n)} = \hat{f}_n. \quad (IV.3.2)$$

Нетрудно видеть, что если нормы $\|A_{nn}^{-1} \Gamma_{nk}\|$ достаточно малы, то операторы \hat{A}_{nn} ограниченно обратимы, а произведения $\hat{A}_{nn} \hat{A}_{nk}$ ограничены, и что если справедливо неравенство (E.7.4), то справедливо и аналогичное неравенство

$$n \geq 0, \|x^{(n)}\| \leq \hat{C}_0 \max_{j \leq n} \|\hat{A}_{jj}^{-1} \hat{f}^{(j)}\|, \quad \hat{C}_0 = \text{const.} \quad (\text{IV.3.3})$$

Точное решение системы (IV.3.2) есть

$$n \geq 0, \quad x^{(n)} = - \sum_{k=0}^{n-1} \hat{A}_{nk}^{-1} \hat{A}_{nk} x^{(k)} + \hat{A}_{nn}^{-1} \hat{f}^{(n)}. \quad (\text{IV.3.4})$$

Вычисление правой части равенства (IV.3.4) обычно связано с некоторой погрешностью. Пусть $U^{(1)}, U^{(2)}, \dots, U^{(n-1)}$ - некоторые точно известные элементы, и пусть $\bar{x}^{(n)}$ - точные значения выражений

$$- \sum_{k=0}^{n-1} \hat{A}_{nk}^{-1} \hat{A}_{nk} U^{(k)} + \hat{A}_{nn}^{-1} \hat{f}^{(n)}; \quad (\text{IV.3.5})$$

примем, что в результате ошибок округления мы, вычисляя величину (IV.3.5), получим не элемент $\bar{x}^{(n)}$, а некоторый другой элемент $U^{(n)} = \bar{x}^{(n)} + d^{(n)}$. Здесь $d^{(n)}$ - результирующая ошибка округления, о которой мы допустим, что она удовлетворяет неравенству

$$\|d^{(n)}\| \leq \varepsilon \|U^{(n)}\|, \quad (\text{IV.3.6})$$

где ε - заранее заданное положительное число, которое может быть сколь угодно малым. Мы принимаем, следовательно, что элемент $\bar{x}^{(n)}$ вычисляется с относительной погрешностью, не превосходящей ε .

2^o. Обозначим $-\hat{A}_{nk}^{-1} \hat{A}_{nk} = a_{nk}$, $\hat{A}_{nn}^{-1} \hat{f}^{(n)} = b^{(n)}$. В этих обозначениях система (IV.3.4) и неравенство (IV.3.3) примут вид

$$n \geq 0, \quad x^{(n)} = \sum_{k=0}^{n-1} a_{nk} x^{(k)} + b^{(n)} \quad (\text{IV.3.7})$$

$$n \geq 0, \quad \|x^{(n)}\| \leq \hat{C}_0 \max_{j \leq n} \|b^{(j)}\|. \quad (\text{IV.3.8})$$

Решая систему (IV.3.7) и учитывая ошибки округления, получим последовательность элементов

$$U^{(0)} = x^{(0)} + d^{(0)} = b^{(0)} + d^{(0)} =: x^{(0)} + r^{(0)},$$

$$U^{(1)} = a_{10} U^{(0)} + b^{(1)} + d^{(1)} = x^{(1)} + a_{10} d^{(0)} + d^{(1)} =: x^{(1)} + r^{(1)}$$

и т.д. Общие формулы имеют вид

$$v^{(n)} = \sum_{k=0}^{n-1} a_{nk} v^{(k)} + b^{(n)} + d^{(n)}, \quad (\text{IV.3.9})$$

$$v^{(n)} = z^{(n)} + y^{(n)}, \quad (\text{IV.3.I})$$

$$y^{(n)} = \sum_{k=0}^{n-1} a_{nk} y^{(k)} + d^{(n)}. \quad (\text{IV.3.II})$$

Система (IV.3.I0) с точностью до обозначений совпадает с системой (IV.3.7); из (IV.3.8) следует тогда

$$\begin{aligned} \|y^{(n)}\| &\leq \hat{C}_0 \max_{j \leq n} \|d^{(j)}\| \leq \hat{C}_0 \varepsilon \max_{j \leq n} \|v^{(j)}\| \leq \\ &\leq \hat{C}_0 \varepsilon [\max_{j \leq n} \|z^{(j)}\| + \max_{j \leq n} \|y^{(j)}\|]. \end{aligned} \quad (\text{IV.3.II})$$

Допустим, что нормы $\|z^{(n)}\|$ ограничены независимо от n ; это будет, например, если выполнены условия теоремы III.7.I, а погрешности Γ_{nk} и $\delta^{(n)}$ достаточно малы. Пусть $\|z^{(n)}\| \leq \hat{C}_1$, тогда

$$\|y^{(n)}\| \leq \hat{C}_0 \hat{C}_1 \varepsilon + \hat{C}_0 \varepsilon \max_{j \leq n} \|y^{(j)}\|. \quad (\text{IV.3.III})$$

Пусть $k \leq n$. Из (IV.3.III) следует:

$$\|y^{(k)}\| \leq \hat{C}_0 \hat{C}_1 \varepsilon + \hat{C}_0 \varepsilon \max_{j \leq k} \|y^{(j)}\| \leq \hat{C}_0 \hat{C}_1 \varepsilon + \hat{C}_0 \varepsilon \max_{j \leq n} \|y^{(n)}\|.$$

Допустим, что $\max_{j \leq n} \|y^{(j)}\|$ достигается при $j = l$. Тогда $\|y^{(l)}\| \leq \hat{C}_0 \hat{C}_1 \varepsilon + \hat{C}_0 \varepsilon \|y^{(l)}\|$ или $\|y^{(l)}\| \leq \hat{C}_0 \hat{C}_1 \varepsilon / (1 - \hat{C}_0 \varepsilon)$

и, тем более,

$$\|y^{(n)}\| \leq \frac{\hat{C}_0 \hat{C}_1 \varepsilon}{1 - \hat{C}_0 \varepsilon}. \quad (\text{IV.3.IV})$$

Таким образом, при сделанных выше предположениях ошибки округления ограничены и стремятся к нулю вместе с ε равномерно по n .

3°. Близкие результаты получаются и в том случае, когда вычисления ведутся с заданной абсолютной погрешностью. По-прежнему считаем, что при вычислении элемента (IV.3.5) получается элемент $v^{(n)} = \sum_{k=0}^{n-1} a_{nk} v^{(k)} + d^{(n)}$, но на этот раз $\|d^{(n)}\| \leq \varepsilon$; число ε по-прежнему считаем сколь угодно малым. Соотношения (IV.3.7) - (IV.3.II) остаются в силе; цепочка неравенств (IV.3.III) упроща-

ется и приводит к нужному результату:

$$\|x^{(n)}\| \leq \hat{C}_0 \max_{j \leq n} \|d^{(j)}\| \leq \hat{C}_0 \varepsilon; \quad (1.3.15)$$

требование ограниченности норм $\|x^{(n)}\|$ в данном случае не необходимо.

Замечание 1. Оценки (1.3.14) и (1.3.15) можно получить, предполагая, что операторы \hat{A}_{nk} удовлетворяют условию (3.7.11).

Замечание 2. Во многих работах трактуются погрешности округления, возникающие при решении линейных алгебраических систем. Отметим некоторые, весьма немногие, из этих работ. Погрешностям округления прямых методов посвящены книги Уилкинсона [1, 2] и Воеводина [1]. В последней книге много внимания уделено вероятностным оценкам. Из работ, в которых изучаются ошибки округления итерационных процессов для линейных алгебраических систем, отметим более старую статью Голуба [1] и недавнюю работу Боняковского [1]. Результаты статьи Голуба мы кратко изложим в § 5 настоящей главы.

§ 4. Погрешность округления метода наискорейшего спуска

1°. Как было сказано в § 8 гл. III, метод наискорейшего спуска сводится к рекуррентному вычислительному процессу; операторы B_n и свободные члены $f^{(n)}$ определяются формулами (1.3.5) и (3.8.2).

Погрешность округления $\gamma^{(n)} = 0^{(n)} - z^{(n)}$ в общем случае определяется рекуррентным соотношением (1.3.11). Для определенности допустим, что вычисления производятся с заданной абсолютной погрешностью, так что $\|d^{(n)}\| < \varepsilon$. Для метода наискорейшего спуска

$$A_{nn} = I; A_{n, n-1} = -B_n = -(I - \sigma_{n-1} A); A_{nk} = 0$$

$$k \leq n-2; \Gamma_{nk} = 0, k \neq n-1; \Gamma_{n, n-1} = -\vartheta'_n A$$

(см. § 8 гл. III). Отсюда

$$\hat{A}_{nn} = I; \hat{A}_{n, n-1} = -I + (\sigma_{n-1} + \vartheta'_n) A; \hat{A}_{nk} = 0,$$

$$k \leq n-2; a_{n, n-1} = I - (\sigma_{n-1} + \vartheta'_n) A; a_{nk} = 0, k \neq n-1.$$

Уравнение (1.7.11) в данном случае имеет вид

$$\forall n \geq 1, x^{(n)} = [I - (\sigma_{n-1} + \vartheta'_n) A] x^{(n-1)} + d^{(n)}.$$

(1.4.1)

Оценим величину $\mu_n := \|A_{n,n-1}\|$. Допустим, что σ_{n-1} при любом n вычисляется с малой относительной погрешностью ϵ' , так что $|\sigma_n| \leq \epsilon' \sigma_{n-1}$. Оператор $A_{n,n-1}$ — ограниченный самосопряженный; его нижняя и верхняя грани суть соответственно

$$1 - (\sigma_{n-1} + \sigma_n) \mu > 1 - \left(\frac{1}{\mu} + \frac{\epsilon'}{\mu}\right) \mu = -\frac{\mu-m}{\mu} - \epsilon' \frac{\mu}{\mu}$$

$$1 - (\sigma_{n-1} + \sigma_n) m \leq 1 - \left(\frac{1}{\mu} - \frac{\epsilon'}{\mu}\right) m = \frac{\mu-m}{\mu} + \epsilon' \frac{m}{\mu}.$$

Следя

$$\mu_n \leq \frac{\mu-m}{m} + \epsilon' \frac{\mu}{m} := \mu. \quad (IV.4.2)$$

Из этого соотношения вытекает неравенство

$$\forall n > 1, \quad |r^{(n)}| \leq \mu |r^{(n-1)}| + \epsilon; \quad (IV.4.3)$$

где отметим, что $r^{(0)} = 0$, так как $r^{(0)}$ есть погрешность точно заданного элемента X_0 . Решение системы неравенств (IV.4.3) есть

$$\forall n > 1, \quad |r^{(n)}| \leq \frac{\mu^n - 1}{\mu - 1} \epsilon. \quad (IV.4.4)$$

Если $\mu < 2m$, то при достаточно малом ϵ' будет $\mu < 1$ и

$$|r^{(n)}| < \frac{\epsilon}{1-\mu}. \quad (IV.4.5)$$

В этом случае погрешность округления ограничена независимо от n и имеет порядок $O(\epsilon)$. Если же $\mu \geq 2m$, то не исключено, что с возрастанием n погрешность округления может неограниченно возрастать. Тот же по существу результат получается, если округление производится с заданной относительной погрешностью (см. статью автора [38]).

2°. Рассмотрим пример. В пространстве $L_2(0, 1)$ рассмотрим уравнение

$$(at+1)x(t) = 1; \quad a = \text{const} > 0 \quad (IV.4.6)$$

с очевидным решением $x = (at+1)^{-1}$. Оператор умножения на функцию $at+1$ — положительно определенный и ограниченный, и уравнение (IV.4.6) можно решить методом наискорейшего спуска; за начальное приближение возьмем $X_0 = 0$. Вычисления проведем для значений $a = 1/2, 1, 3, 10$. В нашем примере $m = 1$, $\frac{\mu-m}{m} = \frac{a}{a+2}$ и

следовательно, $\frac{M-m}{M+m} = \frac{a}{a+2}$ и $\frac{M-m}{m} = a$.

Наименьшее число итераций, необходимое для получения теоретической погрешности, меньшей, чем 10^{-4} , было подсчитано по формуле (Ш.8.3), значения этого числа приведены в таблице 1; результаты вычислений приведены в таблице 2. Они, как и следовало ожидать, вполне удовлетворительны для значения $a = 1/2$, при котором $M < 2m$; они оказываются удовлетворительными также при $a = 1$, когда $M = 2m$. Большие значения a приводят к следующему результату: при очень большом числе итераций точность значений функции $X(b)$ повышается, затем эта точность резко падает; далее наступает переполнение и вычисление последующих итераций становится невозможным. Эти явления находятся в соответствии с результатами данного параграфа и § 8 гл. III.

Таблица 1. Число необходимых итераций

$a = \frac{M-m}{m}$	1/2	1	3	10
$\frac{M-m}{M+m}$	1/5	1/3	1/2	5/6
число необходимых итераций	6	9	19	51

Таблица 2. Точные и приближенные значения искомой функции

из-ды	$a = 1/2$		$a = 1$		$a = 3$		
	точное решение	результат 6 итераций	точное решение	результат 9 итераций	точное решение	результат 10 итераций	результат 15 итераций
0,0	1,0000	0,9920	1,0000	0,9936	1,0000	0,9525	-21609
0,1	0,9524	0,9484	0,9091	0,9081	0,7592	0,7601	-16654
0,2	0,9091	0,9075	0,8333	0,8331	0,6250	0,6240	- 5187
0,3	0,8696	0,8691	0,7692	0,7692	0,5265	0,5254	639,9
0,4	0,8333	0,8333	0,7143	0,7143	0,4545	0,4545	440,7
0,5	0,8000	0,8000	0,6667	0,6667	0,4000	0,4000	33,34
0,6	0,7692	0,7692	0,3250	0,3250	0,3571	0,3571	- 48,32
0,7	0,7407	0,7410	0,5882	0,5882	0,3223	0,3226	- 1792
0,8	0,7143	0,7152	0,5556	0,5556	0,2912	0,2942	17528
0,9	0,6887	0,6919	0,5238	0,5267	0,2703	0,2747	135,69
1,0	0,6667	0,6711	0,5000	0,5010	0,2500	0,2527	455,27

Продолжение таблицы 2

λ-	λ = 10		
	точное решение	результат 8 итераций	результат 10 итераций
0,0	1,0000	0,9238	3983
0,1	0,5000	0,5016	-35,8
0,2	0,3333	0,3333	23,97
0,3	0,2500	0,2556	-79,02
0,4	0,2000	0,2052	-7152
0,5	0,1667	0,1667	-13,63
0,6	0,1429	0,1398	8091
0,7	0,1250	0,1236	4940
0,8	0,1111	0,119	-3715
0,9	0,1000	0,0995	2581
1,0	0,0909	0,0930	-14112

§ 5. Построение округления метода Рундсона для линейных алгебраических систем

В качестве второго примера, к которому могут быть применены преобразования § 3 данной главы, мы коротко изложим здесь содержание части статьи Лолуба [1].

В пространстве R_m рассмотрим векторное уравнение

$$(I - B)x = f. \tag{IV.5.1}$$

Здесь I - единичная матрица в R_m , B - симметричная матрица порядка $m \times m$, спектральный радиус которой $\rho < 1$; следовательно, ограниченный оператор $(I - B)^{-1}$. Примем, что $\text{tr} B$ и вектор f заданы точно, без погрешности. По методу Рундсона решение уравнения (IV.5.1) строится как предел векторов $x^{(n)}$, которые вычисляются по двухшаговому рекуррентному процессу

$$x^{(n)} - \omega Bx^{(n-1)} + (\omega - 1)x^{(n-2)} = \omega f; \tag{IV.5.2}$$

$$\omega = \frac{2}{1 + \sqrt{1 - \rho^2}}; \tag{IV.5.3}$$

начальные приближения $x^{(0)}$ и $x^{(1)}$ заданы произвольно.

Рассмотрим несколько более общий рекуррентный процесс

$$t^{(n)} - \omega B t^{(n-1)} + (\omega - 1) t^{(n-2)} = h^{(n)} \quad (IV.5.4)$$

и докажем, что для этого процесса справедливо соотношение (III.7.4). С этой целью построим решения конечной системы

$$k \leq n, \quad t^{(k)} - \omega B t^{(k-1)} + (\omega - 1) t^{(k-2)} = h^{(k)}; \quad (IV.5.5)$$

$t^{(0)}$ и $t^{(1)}$ считаем заданными произвольно.

Матрицу B можно представить в виде $B = Q \Lambda Q^*$, где $\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)$ есть диагональная матрица собственных чисел матрицы B , Q - ортогональная матрица, составленная из собственных векторов матрицы B , Q^* - матрица, сопряженная к Q . В статье Голуба [1] доказывалось, что решение системы (IV.5.5) можно представить в виде

$$k \leq n, \quad t^{(k)} = \sum_{j=1}^{k-1} Q S_{kj} Q^* h^{(j)}; \quad (IV.5.6)$$

S_j - диагональные матрицы, элементы которых, лежащие на главной диагонали, определяются следующим образом: если

$\theta_j = \arccos(\lambda_j / \rho)$ и $\omega - 1 = \tau^2$, то

$$(S_j)_{jj} = \begin{cases} \tau^{j-1} \frac{\sin v \theta_j}{\sin \theta_j}, & \theta_j \neq 0, \pi; \\ \tau^{j-1}, & \theta_j = 0; \\ (-1)^{j-1} \tau^{j-1}, & \theta_j = \pi. \end{cases}$$

Отсюда

$$\|Q S_{kj} Q^* h^{(j)}\| \leq (\kappa - j) \tau^{k-j-1} \|h^{(j)}\|;$$

по формуле (IV.5.6)

$$k \leq n, \quad \|t^{(k)}\| \leq \max_{j \leq n} \|h^{(j)}\| \sum_{j=1}^k (\kappa - j) \tau^{k-j-1}. \quad (IV.5.7)$$

Так как $0 < \tau < 1$, то сумма в (IV.5.7) ограничена. Но тогда условие (III.7.4) выполнено, и можно оценить погрешность округления для процесса Рундсона (IV.5.2).

Решая бесконечную систему (IV.5.2), мы из-за погрешностей округления получим вместо последовательности $\{x^{(n)}\}$ некоторую другую последовательность $U^{(n)}$, которая точно удовлетворяет другой бесконечной системе

$$v^{(n)} - \omega B v^{(n-1)} + (\omega-1)v^{(n-2)} = \omega f + \delta^{(n)}. \quad (\text{IV.5.8})$$

Допуская, например, что вычисления ведутся с фиксированной абсолютной погрешностью ε , имеем $\|\delta^{(n)}\| \leq \varepsilon$. Из результатов § 3 данной главы вытекает, что

$$\|v^{(n)}\| := \|v^{(n)} - x^{(n)}\| \leq C\varepsilon, \quad C = \text{const}. \quad (\text{IV.5.9})$$

Постоянную C легко определить. Вычитая (IV.5.2) из (IV.5.8) получаем

$$r^{(n)} - \omega B r^{(n-1)} + (\omega-1)r^{(n-2)} = \delta^{(n)}.$$

Сравнивая это с (IV.5.5), находим по неравенству (IV.5.8)

$$\|r^{(n)}\| \leq \max_{j \leq n} \|\delta^{(j)}\| \sum_{j=1}^n (n-j)\tau^{n-j-1} \leq \left[\frac{1-\tau^n}{(1-\tau)^2} - \frac{\tau^{n-1}}{1-\tau} \right] \max_{j \leq n} \|\delta^{(j)}\|. \quad (\text{IV.5.10})$$

Несколько усилив последнее неравенство, получим

$$\|r^{(n)}\| \leq \frac{\varepsilon}{(1-\tau)^2}$$

и, следовательно,

$$C = \frac{1}{(1-\tau)^2}. \quad (\text{IV.5.11})$$

Нетрудно получить аналогичные оценки в предположении, что вычисления выполняются с заданной относительной погрешностью.

Значительная часть статьи Голуба [1] посвящена вероятностным оценкам погрешности округления для метода Рундсона. Мы не будем здесь на этом останавливаться.

ГЛАВА V.

Погрешности метода конечных элементов

§ I. Обзор некоторых результатов (погрешность аппроксимации)

1°. Метод конечных элементов (МКЭ) в настоящее время является одним из самых распространенных и часто применяемых методов приближенного решения дифференциальных уравнений, особенно уравнений эллиптического типа. Этому методу посвящена обширная литература, как монографическая, так и журнальная. Мы отметим здесь только монографии Варга [1], автора этой книги [22], Обана [1], Оганесяна и Руховца [1], Стэнга и Фликса [1, 2], Сьярле [1]. Как кажется автору, в литературе, посвященной теории МКЭ, изучаются в основном две проблемы этой теории: точность МКЭ для различных классов задач механики, математической физики и т.п. и погрешности аппроксимации вычислительных процессов МКЭ (по терминологии настоящей книги). Небольшое количество работ посвящено другим погрешностям МКЭ; так, в некоторых работах автора исследована устойчивость МКЭ относительно искажения (см., например, [22]); в некоторых работах, из которых мы здесь отметим работу Руховца [1], исследовано число обусловленности матрицы МКЭ. В настоящем параграфе мы дадим весьма краткий обзор результатов, относящихся к погрешностям МКЭ. Мы не будем добиваться полноты обзора и остановимся только на результатах, которые кажутся нам наиболее важными.

2°. Для оценки погрешности аппроксимации МКЭ основную роль играет теорема об аппроксимации произвольной функции заданного функционального пространства линейными агрегатами координатных функций МКЭ. Если сетка "регулярная", т.е., любой геометрический элемент получается из некоторой фиксированной геометрической фигуры преобразованием масштаба и сдвигом на целочисленный вектор, то координатные функции МКЭ обычно можно получать такими же преобразованиями независимых переменных из некоторого конечного набора функций Δ , удовлетворяющих определенным требованиям гладкости и имеющих компактный носитель. Мы называем эти функции "исходными".

В общем случае упомянутые выше линейные агрегаты могут не аппроксимировать функции заданного класса. Для того, чтобы такая аппроксимация имела место, необходимо и достаточно, чтобы исходные функции удовлетворяли некоторому конечному числу соот-

ношений. Смысл этих соотношений состоит в том, что некоторые линейные комбинации исходных функций со сдвинутыми аргументами представляют собой полиномы (см. Стрэнг и Бико [1], Стрэнг [1] автор [12, 22], Обен [1]). Поясним это на примерах исходных функций, рассмотренных автором и Обеном.

3°. Пусть t - точка евклидова пространства R_m размерности m . В упомянутых выше работах автора [12, 22] определяется понятие системы исходных функций размерности m . Это набор функций $\omega_{qs}(t)$, где $t \in R_m$, q - мультииндекс размерности m , который принимает конечное число значений, зависящее от натурального параметра δ ; функции ω_{qs} удовлетворяют следующим условиям:

$$\omega_{qs} \in C^{(s-1)}(R_m) \cap W_p^{(s)}(R_m), \text{ где } p \text{ какое-нибудь число из промежутка } 1 < p < \infty; \text{supp } \omega_{qs} \subset \{t \in R_m : 0 \leq t \leq \underline{2}\},$$

где символ \underline{a} означает вектор, все составляющие которого в данной системе координат равны числу a ; наконец, $D^{(a)}\omega_{qs}(t) = \delta_{aq}$; $0 < |a|, |q| \leq s-1$. Координатные функции МКЭ определяются формулой

$$y_{ajh} = \omega_{qs}\left(\frac{t}{h} - j\right); j \in Z^m, h > 0. \quad (\text{У.И.1})$$

Доказывается следующее утверждение: множество линейных агрегатов функций (У.И.1) плотно в $W_p^{(s)}(R_m)$ тогда и только тогда, когда исходные функции удовлетворяют соотношениям, которые автор называет фундаментальными:

$$0 \leq |r| \leq s, \sum_{|q|=0}^{s-1} \sum_{t \in I} \frac{(1-i)^{r-q}}{(r-q)!} \omega_{qs}^-(t+i) = \frac{t^r}{r!}; \quad (\text{У.И.2})$$

$$0 \leq t \leq \underline{1};$$

здесь I - совокупность радиусов-векторов вершин единичного куба $0 \leq t \leq \underline{1}$. Может случиться, что фундаментальные соотношения выполнены для $0 \leq |r| \leq \tilde{s}$ при $\tilde{s} > s$. В этом случае мы говорим об усвоенных фундаментальных соотношениях.

Пусть $\tilde{s} \geq s$, и фундаментальные соотношения верны при $0 \leq |r| \leq \tilde{s}$, и пусть $x(t)$ - произвольная функция класса $W_p^{(s+1)}(R_m)$. Тогда при любом $h > 0$ существует такой линейный агрегат $x^h(t)$ функций (У.И.1), что

$$\forall \tilde{s} \leq s, \|x - x^h\|_{W_p^{(s)}(R_m)} \leq C h^{\tilde{s}-s+1} \max_{|a|=\tilde{s}+1} \|D^a x\|_{L_p(R_m)}. \quad (\text{У.И.3})$$

Заметим, что оценка (У.1.3) остается справедливой, если заменить $W_p^{(j)}$ и $W_p^{(j+1)}$ на $C^{(j)}$ и $C^{(j+1)}$.

Близкие результаты, но только для $p = 2$, содержатся у Фикса и Стрэнга [1] и Стрэнга [1]. Для $\delta = 5$ неравенство (У.3.1) дано в статье Демьяновича и автора [1].

Рассмотрим простейший случай, когда $m = 1$, и пусть исходная система состоит из δ функций $\omega_{q_s}(t)$, $0 \leq q \leq \delta - 1$. Наибольшее значение δ , для которого могут быть выполнены усиленные фундаментальные соотношения, равно $2\delta - 1$. Соответствующие исходные функции определяются однозначно: они равны нулю вне отрезка $0 < t < 2$; на каждом из отрезков $0 < t < 1$ и $1 < t < 2$ они суть полиномы степени $2\delta - 1$. Если положить $\omega_{q_s}(t) = \sigma_{q_s}(t+1)$,

то

$$-1 < t < 0, \sigma_{q_s}(t) = \frac{(1+t)^\delta t^q}{q!} \sum_{k=0}^{\delta-q-1} a_k^{(s)} t^k;$$

$$0 < t < 1, \sigma_{q_s}(t) = \frac{(1-t)^\delta t^q}{q!} \sum_{k=0}^{\delta-q-1} (-1)^k a_k^{(s)} t^k;$$

(У.1.4)

$$a_k^{(s)} = \binom{\delta+k+1}{k}.$$

Можно отметить, что сумма в первой из формул (У.1.4) есть отрезок ряда Тейлора функции $(1-t)^{-\delta}$.

4°. Из неравенства (У.1.3) легко вытекает оценка погрешности аппроксимации МКЭ в энергетической метрике. Рассмотрим эллиптическое уравнение (скалярное или векторное - безразлично) в конечной области $\Omega \subset R_m$ при таких краевых условиях, при которых оператор задачи положительно определен в $L_2(\Omega)$. Относительно области Ω предположим следующее: если все краевые условия - естественные, то Ω может быть произвольной областью с липшицевой границей; если среди краевых условий есть главные, то предположим следующее: существует такая последовательность $h \rightarrow 0$, что при любом h из этой последовательности область Ω может быть исчерпана кубами сетки с длиной ребра $2h$.

Пусть исходная система ω_{q_s} имеет размерность m , и пусть она удовлетворяет усиленным фундаментальным соотношениям при $0 < |t| \leq \delta$, $\delta > 5$. Наконец, допустим, что энергетическая норма

задачи эквивалентна форме $W_2^{(s+1)}$. Тогда

$$|x_* - x_*^{(h)}| \leq C! \bar{\delta}^{-s+1} \max_{|d|=s+1} \|D^d x_*\|_{2, \Omega}; \quad (Y.I.5)$$

здесь x_* - точное решение задачи, $x_*^{(h)}$ - ее приближенное решение по МКЭ.

Если среди краевых условий есть главные, а Ω имеет произвольную Lipschitz-границу, то удается установить значительно худшую оценку: $|x_* - x_*^{(h)}| = O(\sqrt{h})$.

5°. Представляет интерес оценка постоянной C в неравенстве (Y.I.3). Такие оценки получены для ряда случаев в работах автора [24, 31]). Здесь отметим случай, который кажется нам более важным. Пусть $m = 1$, исходные функции определены формулами (Y.I.4) и в формуле (Y.I.3) $W_p^{(3)}(R_m)$ и $W_p^{(s+1)}(R_m)$ заменены на $C^{(\delta)}(R_m)$ и $C^{(s+1)}(R_m)$; пусть еще $\bar{\delta} = 2\delta - 1$. Постоянная C в неравенстве (Y.I.3) зависит от δ и $\bar{\delta}$; обозначим ее через $C(\delta, \bar{\delta})$. Доказывается справедливость оценки

$$C(\delta, \bar{\delta}) < \frac{1}{(2\delta - \bar{\delta})!} + \frac{13e^{3/2} - 2e^{3/2}}{5} \frac{2^{2\delta - 2} \mu(\delta, \bar{\delta})}{(\delta + 1)!} < \\ \leq \frac{1}{(2\delta - \bar{\delta})!} + \frac{8,54 \mu(\delta, \bar{\delta}) \cdot 2^{2\delta - 2}}{(\delta + 1)!}; \quad (Y.I.6)$$

здесь $\mu(\delta, \bar{\delta})$ - функция двух целочисленных аргументов, определяемая формулой

$$\mu(\delta, \bar{\delta}) = \begin{cases} (\delta - 1)^2 (\delta - 2)^2 \dots (\delta - \frac{\bar{\delta} - 1}{2})^2, & \bar{\delta} \text{ - нечетное,} \\ (\delta - 1)^2 (\delta - 2)^2 \dots (\delta - \frac{\bar{\delta}}{2} + 1)^2 (\delta - \frac{\bar{\delta}}{2}), & \bar{\delta} \text{ четное.} \end{cases} \quad (Y.I.7)$$

6°. Обан ([1], гл.V) также рассматривал кубическую сетку в R_m ; его система исходных функций состоит из единственной функции $\mu(t) \in W_2^{(3)}(R_m)$, имеющей компактный носитель. Пусть x произвольная функция класса $W_2^{(s+1)}(R_m)$. Доказывается, что существует функция вида

$$x^h(t) = \sum_{j \in Z^m} a_j \mu(\frac{t}{h} - j), \quad (Y.I.8)$$

удовлетворяющая неравенству

$$\|x - x^h\|_{(3)} \leq Ch^{3+5} \|x\|_{(3+1)}, \quad (V.I.9)$$

если исходная функция $\mu(t)$ удовлетворяет фундаментальным соотношениям, а именно по терминологии Обана, "критерию сходимости")

$$\sum_{k \in \bar{Z}^m} \frac{k^j}{j!} \mu(x-k) = \sum_{0 \leq p \leq j} b^{j-p} \frac{x^p}{p!}; \quad |j| \leq 3, \quad \int_{R_m} \mu(t) dt = 1. \quad (V.I.10)$$

Здесь \bar{Z}^m - множество всех m -мерных целочисленных векторов, b^{j-p} - некоторые постоянные. Условие (V.I.10) не только достаточно, но и необходимо. Отметим, что довольно близкие результаты содержатся в несколько более раннем препринте Фикса и Стрэнга [1].

Приведем еще один пример из книги Обана.

Пусть $\pi_k(t)$ - характеристическая функция интервала $0 < t < 1$, а $\pi_k^{(j)}$, $k = 2, 3, \dots$, последовательные овертки:

$$\pi_k^{(j)}(t) = \int_{-\infty}^{+\infty} \pi_k^{(j-1)}(t-\tau) \pi_k^{(j-1)}(t-\tau) d\tau = \int_0^1 \pi_{k-1}^{(j-1)}(t-\tau) d\tau. \quad (V.I.11)$$

Очевидно, $\pi_{3+2}^{(j)} \in W_9^{(j)}(R_1)$; доказывается, что функция $\pi_{3+2}^{(j)}(x)$ удовлетворяет соотношениям (V.I.10), при этом

$$b^j = \sum_{j_1 + \dots + j_{3+2} = j} \frac{(-1)^j}{(j_1+1)! \dots (j_{3+2}+1)!}. \quad (V.I.12)$$

7⁰. В своей статье [4] Шапошникова рассмотрела среди других вопрос об оценке погрешности производных порядка δ решения первой краевой задачи для уравнения

$$\sum_{|\alpha|, |\beta|=0} (-1)^{|\alpha|} \mathcal{D}^\alpha (A_{\alpha\beta} \mathcal{D}^\beta x) = f(t) \quad (V.I.13)$$

в m -мерном кубе Ω . Как обычно, предполагается, что оператор задачи положительно определенный в $L_2(\Omega)$, уравнение (V.I.13) в Ω не вырождается и коэффициенты $A_{\alpha\beta}$ - достаточно гладкие.

Допустим, что точное решение задачи $x_* \in W_\rho^{(l)}(\Omega)$, где $l < \rho < \infty$ и $l > \delta$; через $x_*^{(h)}(t)$ обозначим приближенное решение, полученное по МКЭ. Доказывается, что справедлива следующая оценка погрешности аппроксимации:

$$|x_* - x_*^{(h)}|_{W_p^{(\nu)}(\Omega)} \leq ch^{\nu - s + \frac{m}{p} - \frac{m}{p}} \omega_{L_p}(x_*, h), \quad \rho \geq 2;$$

$$|x_* - x_*^{(h)}|_{W_p^{(\nu)}(\Omega)} \leq ch^{\nu - s + \frac{m}{p} - \frac{m}{p}} \omega_{L_p}(x_*, h), \quad 1 < \rho < 2. \quad (Y.I.14)$$

Здесь $\omega_{L_p}^{(\nu)}(x_*, h)$ - наибольший L_p -модуль непрерывности произвольных порядка ν решения $x_*(t)$.

Джинкаррани [13] получили оценку

$$|x_* - x_*^{(h)}|_{(3)} \leq ch^{\nu-3} |x_*|_{(4)} \quad (Y.I.15)$$

в предположении, что $\nu > \tilde{\nu} > 3$ и что краевые условия могут быть произвольными - требуется только, чтобы оператор данной задачи был локально определенным в $L_2(\Omega)$.

8°. Небольшое число работ посвящено оценке погрешности аппроксимации МКЭ для вырождающихся дифференциальных уравнений. Не приводя результатов, отметим работу Гусмана и Оганесяна [1], в которой рассмотрено слабо вырождающееся алгебраическое уравнение в прямоугольнике на двумерной плоскости, и работу автора [14] (см. также [22]), в которой рассмотрено линейное обыкновенное дифференциальное уравнение второго порядка с любым показателем вырождения.

§ 2. Обзор некоторых результатов; погрешность искажения и число обусловленности МКЭ.

1°. Сформулируем полученные нами (см. работы автора [12, 13, 22]) результаты, которые относятся к устойчивости МКЭ относительно искажения и к оценке числа обусловленности матрицы МКЭ; о значении числа обусловленности для оценок погрешностей сказано выше, в гл. IV.

Рассмотрим алгебраическое дифференциальное уравнение

$$\sum_{|\alpha|, |\beta|=0}^d (-1)^{|\alpha|} D^\alpha (A_{\alpha\beta} D^\beta x) = f(t); \quad (Y.2.1)$$

будем считать, что t меняется в n -ичной области $\Omega \subset R_m$ с внешней границей. Краевые условия будем считать такими, что оператор задачи (Y.2.1) при этих условиях - положительно определенный в $L_2(\Omega)$. Примем еще, что $f \in L_2(\Omega)$ и что дифференциальный оператор (Y.2.1) не вырождается в $\bar{\Omega}$. Наконец, примем, что $\mathcal{K}(t)$ есть k -мерный вектор, $k \geq 1$.

Нашу краевую задачу будем приближенно решать по МКЭ при координатных функциях, описанных в п. 3^о § I. Устойчивость и число обусловленности последуем в таких предположениях: либо речь идет о первой краевой задаче, тогда область Ω может быть произвольной липшицевой, либо Ω - куб, тогда краевые условия могут быть произвольными.

Рассмотрим сначала первую краевую задачу. Пусть h - шаг сетки. Рассмотрим открытые кубы сетки с ребром h . Радиус-векторы их верш имеют вид j^h , где $j \in \mathbb{Z}^m$ есть целочисленный вектор, который будем называть номером соответствующей вершины. У каждого куба сетки есть вершина с наименьшим номером j . Обозначим через \mathcal{J}^h множество таких наименьших номеров, соответствующих кубам сетки, лежащим внутри Ω . Приближенное решение задачи будем искать в виде

$$\alpha^h(t) = \sum_j \sum_{j \in \mathcal{J}^h} a_{qj}^{(h)} \varphi_{qjh}(t), \quad (У.2.2)$$

где φ_{qjh} - функции (У.1.1). Обозначим через $a^{(h)}$ вектор коэффициентов $a_{qj}^{(h)}$ и через N его размерность. Этот вектор определяется из алгебраической системы

$$A_h a^{(h)} = f^{(h)}, \quad (У.2.3)$$

где A_h - матрица порядка $N \times N$, составленная из скалярных произведений $[\varphi_{qjh}, \varphi_{q'j'h}]$, а $f^{(h)}$ - вектор с составляющими (f, φ_{qj^h}) .

Наряду с обычным евклидовым пространством N измерений R_N рассмотрим два N -мерных гильбертовых пространства \bar{X}_h и \bar{Y}_h с нормами

$$\|a\|_{\bar{X}_h} = h^{m/2} \|a\|_{R_N}, \quad \|a\|_{\bar{Y}_h} = h^{-m/2} \|a\|_{R_N}. \quad (У.2.4)$$

Пусть $\{h_k\}$ - последовательность значений h , стремящаяся к нулю, и $\{N_k\}$ соответствующая последовательность размерностей векторов коэффициент α . Доказывается, что процесс (У.2.3) определения вектора коэффициентов $\alpha^{(h)}$ устойчив (в смысле второго определения) в последовательности пар пространств $(\bar{X}_{h_k}, \bar{Y}_{h_k})$. Существуют, следовательно, такие постоянные $\rho, q, \nu > 0$, что если $\|\bar{f}_{h_k}\| \leq \nu$, то искаженное уравнение

$$(A_{h_k} + \Gamma_{h_k}) b^{(h_k)} = f^{(h_k)} + \delta^{(h_k)}$$

имеет единственное решение и

$$\| \bar{v}^{(h_k)} - a^{(h_k)} \|_{\bar{x}_{h_k}} \leq \rho \| \Gamma_{h_k} \|_{\bar{x}_{h_k} \rightarrow \bar{y}_{h_k}} + q \| \delta^{(h_k)} \|_{\bar{y}_{h_k}}. \quad (V.2.5)$$

Далее, пусть

$$\bar{z}^{(h_k)}(t) = \sum_q \sum_{j \in J_{h_k}} v_{qj}^{(h_k)} y_{qj|_{h_k}}$$

- покаженное приближенное решение. Доказывается, что существуют такие постоянные $\rho', q', \nu' > 0$, что если $\| \Gamma_{h_k} \| \leq \nu'$, то

$$\| \bar{z}^{(h_k)} - x^{(h_k)} \| \leq \rho' \| \Gamma_{h_k} \| + q' \| \delta^{(h_k)} \|; \quad (V.2.6)$$

иначе говоря, доказывается устойчивость процесса вычисления приближенного решения в соответствующей последовательности пар пространств.

Для той же первой краевой задачи доказывается, что число обусловленности матрицы A_{h_k} имеет порядок $O(h^{-23})$; см. работы автора [1, 22].

Те же результаты получаются, если область Ω - куб, а краевые условия - естественные. Эти результаты верны и для процесса вычисления собственных чисел по МКЭ; см. книгу автора [22].

2°. Руховец [1] предложил некоторое простое видоизменение функционала энергии для уравнения (V.2.1) и соответствующих естественных краевых условий, в результате которого точное решение задачи не меняется, а оценка $O(h^{-23})$ для числа обусловленности матрицы МКЭ справедлива в случае любой области с достаточно гладкой границей. Как доказал Чистов [1], для функционала Руховца справедливы приведенные в § I данной главы теоремы об устойчивости процесса МКЭ, если граница области Ω достаточно гладкая.

3°. Первой работой, в которой исследовалась устойчивость процесса МКЭ, была статья Демьяновича [1]. В этой статье рассмотрена несамосопряженная эллиптическая задача

$$-\sum_{i,k=1}^m \frac{\partial}{\partial t_i} (a_{ik} \frac{\partial x}{\partial t_k}) + \sum_{i=1}^m a_i \frac{\partial x}{\partial t_i} + ax = f(t), \quad t \in \Omega. \quad (V.2.7)$$

На область Ω накладываются некоторые ограничения, которые выполняются, например, если $\partial \Omega \in C^{(2)}$. Вводится кусочно полиномиальная походная функция; с ее помощью строятся координатные функции МКЭ. Вводится в рассмотрение пространство сеточных функ-

ций с узлами, лежащими внутри Ω ; на этом пространстве вводят-ся две гильбертовы метрики, $(\cdot, \cdot)_0$ и $(\cdot, \cdot)_1$, которые в пределе, при $h \rightarrow 0$, переходят соответственно в метрику L_2 и в энергетическую метрику главной части оператора задачи. Норма произвольной матрицы A определяется формулой

$$\|A\| = \sup \frac{|(Av, v')_0|}{\|v\|_1 \cdot \|v'\|_1}; \quad (У.2.8)$$

здесь v, v' — произвольные элементы упомянутого выше пространства сеточных функций. Пусть $\Gamma_h, \delta^{(h)}, z^{(h)} - x^{(h)}$ суть соответственно искажения матрицы МКЭ, столбца свободных членов и приближенного решения; само приближенное решение обозначим через $x^{(h)}$. Доказывается существование таких постоянных $p, q, \epsilon > 0$, что если

$$\|\Gamma_h\| < \epsilon, \quad \|z^{(h)} - x^{(h)}\|_1 < p \|\Gamma_h\| + q \|\delta^{(h)}\|_1. \quad (У.2.9)$$

В той же статье исследуется устойчивость процесса МКЭ для уравнения (У.2.7) в предположении, что искажение Γ_h принадлежит некоторому специальному пространству матриц. Идеи работы Демьяновича [1] развиты в последующих работах того же автора [2, 3].

§ 3. МКЭ с одной кусочно полиномиальной исходной функцией

1°. В ближайших параграфах данной главы мы подробнее исследуем погрешности МКЭ в том случае, когда исходная система содержит только одну функцию, удовлетворяющую фундаментальным соотношениям (У.1.10). Для простоты ограничимся более частным случаем, когда размерность $m = 1$, а исходная функция есть кусочно полиномиальная функция $\pi_\kappa(t)$ (формула (У.1.11)). Погрешность аппроксимации в этом случае исследована в работах Стрэнга и Фикса [1] и Осэна [1]; см. также п. 6° § 1 данной главы, и в этом плане мы ограничимся тем, что оценим постоянную C в неравенстве (У.1.9). Этому будет посвящен настоящий параграф. В § 4 мы рассмотрим погрешность искажения, а в § 5 — число обусловленности матрицы МКЭ для той же исходной функции $\pi_\kappa(t)$. В § 6 мы коротко рассмотрим погрешности МКЭ для симплициальных сеток.

2°. Начнем с вывода оценки производных полиномов $\pi_\kappa(t)$. Имеем

$$\pi_0(t) = \begin{cases} 1, & t \in (0, 1), \\ 0, & t \notin (0, 1) \end{cases}$$

и

$$\pi_k(t) = \int_0^1 \pi_{k-1}(t-\tau) d\tau; \quad (У.3.1)$$

отсюда, очевидно,

$$0 \leq \pi_{k-1}(t) \leq 1. \quad (У.3.2)$$

Замена $t-\xi = \tau$ дает

$$\pi_k(t) = \int_{t-1}^t \pi_{k-1}(\xi) d\xi$$

и, следовательно,

$$k \geq 1, \pi_k^{(\delta)}(t) = \pi_{k-1}^{(\delta-1)}(t) - \pi_{k-1}^{(\delta-1)}(t-1). \quad (У.3.3)$$

В частности, при $\delta = 1$, $|\pi_k'(t)| \leq \pi_{k-1}(t) + \pi_{k-1}(t-1) \leq 2$, и индукция

$$|\pi_k^{(\delta)}(t)| \leq 2^\delta. \quad (У.3.4)$$

С целью оценить постоянную в формуле (У.1.9), мы воспользуемся здесь выводом этой формулы, данный в книге Обэна [1], до тех пор пока оценка констант, вводимых в ходе рассуждений. Начиная с построения аппроксимирующей функции $\mathcal{X}^{(h)}(t)$. Наряду с новой функцией $\mu(t) = \pi_k(t)$ введем в рассмотрение еще одну функцию $\lambda(t)$ со следующими свойствами: носитель этой функции компактен, она ограничена и измерима и

$$\int_a^{\infty} \lambda(t) dt = 1, \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mu(t) \lambda(\tau) (t-\tau)^k dt d\tau = \begin{cases} 1, & k=0 \\ 0, & 0 < k \leq \delta \end{cases}. \quad (У.3.5)$$

Пусть $x \in W_2^{(\delta+1)}(R_1)$, причем $\text{supp } x$ компактен. На каждом интервале $(jh, (j+1)h)$ разложим $x(t)$ по формуле Тейлора в круг точки jh и введем обозначения

$$X_{jh}^\delta(x)(t) = \sum_{q=0}^{\delta-1} \frac{h^q}{q!} x^{(q)}(jh) \left(\frac{t}{h} - j\right)^q,$$

$$Y_{jh}^\delta(x)(t) = \int_{jh}^t \frac{(t-\tau)^{\delta-1}}{(\delta-1)!} x(\tau) d\tau, \quad (У.3.6)$$

тогда $x(t) = X_{jh}^\delta(x)(t) + Y_{jh}^\delta(x)(t)$.

Вводится оператор, который мы обозначим через g_h и который действует по правилу

$$(g_h x)(t) = \sum_{\rho \in \mathbb{Z}} x_h^\rho \mu\left(\frac{t}{h} - \rho\right); \quad \mu(t) = \pi_k(t);$$

$$x_h^\rho = \frac{1}{h} \int_{-\infty}^{+\infty} x(t) \lambda\left(\frac{t}{h} - \rho\right) dt.$$

Доводяется, что этот оператор оставляет инвариантными усеченно полиномиальные функции отепени l .

Положим теперь $x^{(h)}(t) = (g_h x)(t)$ и оценим разность

$$\| \mathcal{D}^j(x(t) - x^{(h)}(t)) \|_{L_2}^2 = \sum_{i \in \mathbb{Z}} U_{ih}^2; \quad \mathcal{D}^j = \frac{d^j}{dt^j};$$

$$U_{ih}^2 = \int_{ih}^{(i+1)h} [\mathcal{D}^j(I - g_h)(Y_{ih}^s(x^{(h)}(t)))]^2 dt \leq \quad (У.3.7)$$

$$\leq 2 \int_{ih}^{(i+1)h} [Y_{ih}^{s-j}(x^{(h)}(t))]^2 dt + 2 \int_{ih}^{(i+1)h} [\mathcal{D}^j g_h(Y_{ih}^s(x^{(h)}(t)))]^2 dt.$$

Первый член справа оценим по неравенству Коши - Буняковского:

$$\int_{ih}^{(i+1)h} [Y_{ih}^{s-j}(x^{(h)}(t))]^2 dt \leq$$

$$\leq \frac{h^{2(s-j)}}{[(s-j-1)!]^2 (2s-2j-1)(2s-2j)} \int_{ih}^{(i+1)h} [x^{(h)}(t)]^2 dt. \quad (У.3.8)$$

Займемся вторым интегралом. Временно положим $Y_{ih}^s(x^{(h)}(t)) = v(t)$.

По определению оператора g_h имеем

$$\mathcal{D}^j g_h(Y_{ih}^s(x^{(h)}(t))) = (\mathcal{D}^j g_h v)(t) = h^{-j} \sum_{\rho \in \mathbb{Z}} v_h^\rho \mu^{(j)}\left(\frac{t}{h} - \rho\right),$$

$$v_h^\rho = \frac{1}{h} \int_{(a+\rho)h}^{(b+\rho)h} \lambda\left(\frac{\tau}{h} - \rho\right) v(\tau) d\tau,$$

где $(a, b) \supset \text{supp } \lambda(t)$. Далее,

$$(\mathcal{D}^j g_h v)^2 \leq h^{-2j} \sum_{\rho \in \mathbb{Z}} (v_h^\rho)^2 \sum_{\rho \in \mathbb{Z}} [\pi_k^{(j)}\left(\frac{t}{h} - \rho\right)]^2. \quad (У.3.9)$$

Каждая из производных $\pi_k^{(j)}\left(\frac{t}{h} - \rho\right)$, отличных от нуля (их число не превосходит $k+1$), не превосходит 2^j , поэтому

$$(\mathcal{D}^j g_h v)^2 \leq (k+1) 2^{2j} h^{-2j} \sum_{\rho \in \mathbb{Z}} (v_h^\rho)^2. \quad (У.3.10)$$

Суммировать в (У.3.10) достаточно по тем ρ , для которых $\rho \leq$

$\frac{1}{h} - \rho \leq k+1$; число τ из ρ не превосходит $k+1$
 Пусть $|\lambda(t)| \leq C_0 = \text{const}$ Тогда

$$(v_{ih}^\rho)^2 \leq \frac{1}{h} C_0^2 (b-a) \int_{(a+\rho)h}^{(b+\rho)h} v^2(t) dt.$$

Интегрируя неравенство (У.3.10) по интервалу $(ih, (i+1)h)$, получим

$$\int_{ih}^{(i+1)h} (Dq_h v)^2 dt \leq C_0^2 (b-a)^2 h^{2j-2j} \sum \int_{(a+\rho)h}^{(b+\rho)h} v^2(t) dt. \quad (\text{У.3.11})$$

Найдем пределы изменения ρ , когда $ih \leq t \leq (i+1)h$. Число ρ должно удовлетворять неравенствам $i-\rho \geq 0$, $i+1-\rho \leq k+1$, которые означают, что $(\frac{1}{h} - \rho) \in [0, k+1] = \text{supp } \pi_k(t)$. Отсюда следует, что в (У.3.11) $\rho = i-k, i-k+1, \dots, i$, и это неравенство можно заменить следующим:

$$\int_{ih}^{(i+1)h} (Dq_h v)^2 dt \leq C_0^2 (b-a) (k+1)^2 h^{2j-2j} \int_{(a+i-k)h}^{(b+i)h} v^2(t) dt.$$

Вспомним, что $v = Y_{ih}^j(x^{(3)})$. Это приведет нас к неравенству

$$\int_{ih}^{(i+1)h} (Dq_h Y_{ih}^j(x^{(3)}))^2 dt \leq C_0^2 (b-a) (k+1)^2 h^{2j-2j} \int_{(a+i-k)h}^{(b+i)h} (Y_{ih}^j(x^{(3)}))^2 dt.$$

Аналогично неравенству (У.3.8) нетрудно получить неравенство

$$\int_{(a+i-k)h}^{(b+i)h} (Dq_h Y_{ih}^j(x)) ^2 dt \leq \frac{h^{2(s-j)^2} C_0^2 (b-a)^2 (b^{2s} - (a-k)^{2s})}{2s(2s-1)((s-1)!)^2} \int_{ih}^{(i+1)h} (x^{(3)}(t))^2 dt.$$

Отсюда и из (У.3.9) найдем

$$U_{ih}^2 \leq C_1 h^{2(s-j)^2} \int_{ih}^{(i+1)h} (x^{(3)})^2 dt, \quad (\text{У.3.12})$$

где

$$C_1^2 = [(s-j-1)!]^2 [(s-j)^2 2s-2j-1]^{-2} + C_0^2 (b-a) (k+1)^2 \frac{(b^{2s} - (a-k)^{2s})}{s(2s-1)((s-1)!)^2}. \quad (\text{У.3.13})$$

Наконец, просуммировав неравенство (У.3.12) по i , и учтя

$$\int_{-\infty}^{+\infty} [\mathcal{D}(x(t) - x^{(h)}(t))]^2 dt \leq C_1^2 h^{2(s-j)} \int_{-\infty}^{+\infty} (x^{(s)})^2 dt; \quad (\text{У.3.14})$$

величина C_1 определяется формулой (У.3.13).

Допустим, что энергетическая норма нашей задачи имеет оценку

$$|x| \leq C_2 \sum_{j=0}^K \|x^{(j)}\|^2. \quad (\text{У.3.15})$$

Обозначим через x_* точное решение задачи, а через $x_*^{(h)}$ приближенное решение по МКЭ; допустим еще, что $x_* \in W_2^{(s)}$. Тогда, полагая $h < 1$, находим

$$|x_* - x_*^{(h)}| \leq C_3 \sum_{j=0}^j \| \mathcal{D}(x_* - x_*^{(h)}) \|^2 \leq \quad (\text{У.3.16})$$

$$\leq (C_1 C_2)^2 \frac{h^{2(s-K)}}{1-h^2} \|x_*^{(s)}\|^2;$$

постоянные C_1 и C_2 определяются формулами (У.3.13) и (У.3.15).

§ 4. Точность кошения

Г°. Рассмотрим обыкновенное дифференциальное уравнение

$$\sum_{k=0}^s (-1)^k \frac{d^k}{dt^k} [P_k(t) \frac{d^k x}{dt^k}] = f(t), \quad 0 < t < 1. \quad (\text{У.4.1})$$

Коэффициенты $P_k(t)$ предполагаем измеримыми и ограниченными, причем коэффициент $P_s(t)$ положительно ограничен снизу; крайние условия таковы, что оператор задачи положительно определен в $L_2(0, 1)$.

Приближенное решение задачи будем искать в виде

$$x^{(h)}(t) = \sum_{j \in \mathbb{Z}} a_j^{(h)} \pi_j\left(\frac{t}{h} - j\right), \quad h = \frac{1}{(s+1)n},$$

где n - произвольно выбранное натуральное число. Нам интересуют только значения $t \in (0, 1)$, и достаточно суммировать лишь по тем j , для которых

$$\frac{t}{h} - j \in \text{supp } \pi_j = [0, s+1]. \quad (\text{У.4.2})$$

Пусть $t/h = j_0 + \tau$, где j_0 - целое число и $0 \leq \tau < 1$. Тогда $t/h - j = j_0 - j + \tau$, и условие (У.4.2) будет выполнено, если $0 \leq j_0 - j \leq s$, или $j_0 - s \leq j \leq j_0$. Далее, j_0 меняется в пределах

$0 \leq j_0 \leq (s+1)n - 1$, откуда получаются пределы изменения j :
 $-s \leq j \leq (s+1)n - 1$, и мы приходим к следующему выражению

для $x^{(h)}$:

$$x = \sum_{j=-3}^{(h) \text{ mod } 3} a_j \pi_3\left(\frac{t}{h} - j\right). \tag{У.4.3}$$

Система уравнений МКЭ имеет вид:

$$\sum_{j=-3}^{(h) \text{ mod } 3} a_j [y_{3jh}, y_{3jh}] = (f, y_{3jh}); \quad y_{3jh} = \pi_3\left(\frac{t}{h} - j\right); \tag{У.4.4}$$

как обычно, здесь квадратными скобками обозначено скалярное произведение в энергетической метрике, а круглыми скобками - скалярное произведение в $L_2(0, 1)$.

Обозначим через $N = n(\delta+1) + \delta$ порядок системы (У.4.4), через M_h - матрицу этой системы, через $a^{(h)}$ и $f^{(h)}$ - векторы в \mathbb{R}^N составляющими (a_j, y_{3jh}) и (f, y_{3jh}) соответственно. Система (У.4.4) записывается в виде

$$M_h a^{(h)} = f^{(h)}. \tag{У.4.5}$$

2°. Нормы $\|x^{(h)}\|_2$ и $\|a^{(h)}\|_2$ связаны двусторонней оценкой; выведем ее. Имеем

$$\|x^{(h)}\|_2^2 = \sum_{l=1}^{(h) \text{ mod } 3} \int_{(l-1)h}^{lh} [x^{(h)}]^2 dt.$$

На промежутке $((l-1)h, lh)$ отличны от тождественного нуля только функции $y_{3jh}(t)$, $j = l-3, l-2, \dots, l-1$, поэтому

$$\int_{(l-1)h}^{lh} [x^{(h)}]^2 dt = \int_{(l-1)h}^{lh} \left[\sum_{j=l-3}^{l-1} a_j \pi_3\left(\frac{t}{h} - j\right) \right]^2 dt = h \int_0^{\delta} \left[\sum_{p=0}^{\delta} a_{l-3+p} \pi_3\left(\frac{x+p}{3}\right) \right]^2 dx. \tag{У.4.6}$$

Последний интеграл есть неотрицательная квадратичная форма от переменных $a_{l-3}, a_{l-2}, \dots, a_{l-1}$. Докажем, что эта форма положительно определенная. Пусть она равна нулю. Тогда при $0 \leq x \leq 1$

$$\sum_{p=0}^{\delta} a_{l-3+p} \pi_3(x+p) = 0.$$

Легко доказать, что коэффициенты a_{l-3+p} равны нулю, т.е. что функции $\pi_3(x+p)$, $p = 0, 1, 2, \dots, \delta$ линейно независимы. Это утверждение очевидно справедливо при $\delta = 0$ и $\delta = 1$. Допустим, что оно справедливо для некоторого значения $\delta - 1$ и л

кажем, что оно справедливо и для индекса δ . Итак, пусть

$$\sum_{\rho=0}^{\delta} d_{\rho} \pi_{\delta}(\alpha + \rho) = 0, \quad 0 \leq \alpha \leq 1.$$

По определению функций π_{δ}

$$0 = \sum_{\rho=0}^{\delta} d_{\rho} \int_{-\infty}^{+\infty} \pi_{\delta-1}(y) \pi_{\delta}(\alpha + \rho - y) dy = \sum_{\rho=0}^{\delta} d_{\rho} \int_{\alpha+\rho-1}^{\alpha+\rho} \pi_{\delta-1}(y) dy.$$

Продифференцируем по α :

$$\sum_{\rho=0}^{\delta} \beta_{\rho} \pi_{\delta-1}(\alpha + \rho) = 0,$$

где $\beta_0 = -d_0$; $\beta_{\delta} = d_{\delta}$; $\beta_{\rho} = d_{\rho} - d_{\rho+1}$, $0 \leq \rho \leq \delta-1$.

Так как α меняется в пределах $(0, 1)$, то $\pi_{\delta-1}(\alpha-1) = 0$, $\pi_{\delta-1}(\alpha+\delta) = 0$, а члены с индексами $\rho = -1$ и $\rho = \delta$ исчезают. Мы приходим к тождеству

$$\sum_{\rho=0}^{\delta-1} \beta_{\rho} \pi_{\delta-1}(\alpha + \rho) = 0.$$

По индукционному предположению, $\beta_{\rho} = 0$, $0 \leq \rho \leq \delta-1$, или $d_0 = d_1 = \dots = d_{\delta} := d$. Теперь

$$d \sum_{\rho=0}^{\delta} \pi_{\delta}(\alpha + \rho) = 0;$$

так как $\pi_{\delta}(\alpha + \rho) \neq 0$ и $\pi_{\delta}(\alpha + \rho) \geq 0$, то $d = 0$, и наше утверждение доказано. Из этого утверждения вытекает существование положительных и не зависящих от ℓ и h величин C_1 и C_2 , таких, что справедливы неравенства

$$\frac{C_1^2}{\delta+1} \sum_{j=\ell-\delta+1}^{\ell-1} (a_j^{(h)})^2 \leq \int_0^1 \left[\sum_{j=\ell-\delta+1}^{\ell-1} a_j^{(h)} \pi_{\delta} \left(\frac{t}{h} - j \right) \right]^2 dt \leq \frac{C_2^2}{\delta+1} \sum_{j=\ell-\delta+1}^{\ell-1} (a_j^{(h)})^2,$$

или, если воспользоваться формулой (У.4.6),

$$\frac{C_1^2 h}{\delta+1} \sum_{j=\ell-\delta+1}^{\ell-1} (a_j^{(h)})^2 \leq \int_{(\ell-1)h}^{\ell h} [x^{(h)}(t)]^2 dt \leq \frac{C_2^2 h}{\delta+1} \sum_{j=\ell-\delta+1}^{\ell-1} (a_j^{(h)})^2.$$

Просуммировав это по ℓ , получим искомого неравенство

$$C_1 \sqrt{h} \|a\|_{R_r}^{(h)} \leq \|x\|_2 \leq C_2 \sqrt{h} \|a\|_{R_r}^{(h)}. \quad (\text{У.4.7})$$

3° . Введем в рассмотрение N -мерные гильбертовы пространства \bar{X}_r и \bar{Y}_r с нормами

$$\forall a \in R_r, \|a\|_{\bar{X}_r} = \sqrt{h} \|a\|_{R_r}, \quad \|a\|_{\bar{Y}_r} = \sqrt{h^{-1}} \|a\|_{R_r}.$$

эти пространства — сопряженные относительно скалярного умножения в $R_{\mathcal{H}}$.

Обозначим через A_h оператор, действующий из \bar{X}_h в \bar{Y}_h и рожденный матрицей M_h . Будем рассматривать $a^{(h)}$ и $f^{(h)}$ как элементы пространств \bar{X}_h и \bar{Y}_h соответственно. Тогда уравнению (У.4.5) можно придать вид

$$A_h a^{(h)} = f^{(h)}. \quad (\text{У.4.8})$$

Теорема У.4.1. Вычислительный процесс (У.4.8) устойчив в смысле второго определения в последовательности пар пространств (\bar{X}_h, \bar{Y}_h) .

Доказательство этой теоремы дословно совпадает с доказательством теоремы I § 2 главы IX книги автора [22], и мы его здесь приводить не будем.

Обозначим через $H^{(h)}$ подпространство энергетического пространства задачи настоящего параграфа, образованное функциями вида

$$x^{(h)}(t) = \sum_{j=-3}^{n(s+1)-1} d_j \pi_3\left(\frac{t}{h} - j\right) \quad (\text{У.4.9})$$

с произвольными коэффициентами d_j . Оператор (У.4.9), переводя вектор $d = (d_{-3}, d_{-3+1}, \dots, d_{n(s+1)-1})$ в функцию $x^{(h)} \in H^{(h)}$, обозначим через Π_h . Приближенное решение $x^{(h)}(t)$ можно определить из условия

$$A_h \Pi_h^{-1} x^{(h)} = f^{(h)}. \quad (\text{У.4.10})$$

Теорема У.4.2. Вычислительный процесс (У.4.10) устойчив в смысле второго определения в последовательности пар пространств $(H^{(h)}, R_{\mathcal{H}})$.

Эту теорему мы здесь также не будем доказывать: ее доказательство совпадает с доказательством теоремы 2 § 2 главы IX книги автора [22].

Так же, в цитированной здесь книге [22], доказывается, что число обусловленности матрицы M_h (уравнение (У.4.5)) имеет порядок $O(h^{-2s})$.

§ 5. Симплициальные сетки; оценка аппроксимации.

1^o. Симплициальные сетки имеют то важное преимущество перед кубическими, что они, вообще говоря, позволяют лучше аппроксимировать произвольно заданную область. Из большого количества работ, посвященных методу конечных элементов с симплициальными

отками, мы отметим здесь лишь немногие из тех, в которых изучается погрешность аппроксимации или искажения; та же только, что симплицальные сетки систематически используются в книгах Ганесяна и Руховца [1], Стрэнга и Фикса [2] и Сьярле [1]. Краткий обзор работ по МКЭ на симплицальных сетках за время от 1973 г. дан в статье Корнеева [1].

Зламал [1] построил приближенное решение двумерного эллиптического уравнения четвертого порядка в виде функции, которая на каждом треугольнике сетки представляет собой полином 5-й степени. Метод и результаты этой работы были обобщены Женишеком [1], который рассмотрел кусочно полиномиальные функции степени $\delta + 1$ и обосновал метод для значений $\delta = 1, 2$.

Дальнейшее обобщение дано в работах Брамбла и Зламала [1], в которой рассматриваются эллиптические уравнения любого порядка в двумерной области. В каждом из треугольников сетки выделяется совокупность узловых точек, одна из которых лежит в центре тяжести треугольника, а остальные выбираются в вершинах и, по определенному правилу, на сторонах треугольника. В узлах заданы значения полинома и некоторых его производных; максимальные порядки этих производных определяются так, чтобы число условий равнялось $(2\delta + 1)(4\delta + 3)$ - числу коэффициентов полинома степени $\delta + 1$ от двух переменных. Доказываются следующие утверждения:

1) Существует единственный полином $P_\delta(t_1, t_2)$ степени $4\delta + 1$, удовлетворяющий упомянутым выше условиям. Заметим, что в работе Корнеева и Пономарева [1] дан явный вид полиномов $P_\delta(t_1, t_2)$.

2) Пусть Ω - замкнутая многоугольная область, разбитая на замкнутые же треугольники T_τ . Допустим, что узлы, расположенные на сторонах или в вершинах треугольников, выбраны одинаковыми для всех тех треугольников, которым рассматриваемый узел принадлежит. Тогда функция $\psi_\delta(t_1, t_2)$, которая на каждом из треугольников T_τ с задает с соответствующим полиномом $P_\delta(t_1, t_2)$ принадлежит при подходящем выборе δ к классу $C^{(m)}$, где m - половина порядка данного дифференциального уравнения.

3) Пусть дана функция $X \in W_2^{(k)}(\Gamma)$, где $2\delta + 2 \leq k \leq 4\delta + 2$, а Γ - некоторый треугольник с выбранными на нем узлами. Заданы на этих узлах значения функции X и ее соответствующих производных и построим по этим данным полином $P_\delta(t_1, t_2)$ так, как это

указано выше. Тогда справедлива оценка

$$\|x - p_3\|_{(0)} \leq \frac{Ch^{k-l}}{(\sin d)^{5l}} \|x\|_{(k)}, \quad l \leq k-1. \quad (У.5.1)$$

Здесь d - наименьший угол и h - наибольшая сторона треугольника T . Оценка (У.5.1) установлена для всех $k \geq 2$, кроме $k = 3$ для $k = 3$ оценка (У.5.1) получена Чистовым [1].

2°. В своих статьях [2,3] Женишек рассматривает тетраэдральную сетку в трехмерном пространстве. По аналогии с результатами, полученными для треугольной сетки на плоскости, Женишек делает следующее предположение: простейший полином в m -мерном симплексе, который позволяет строить кусочно полиномиальные раз непрерывно дифференцируемые функции, имеет степень $n = 2m$. В цитированных работах Женишека излагаются результаты, полученные им для $m = 3$ и $j = 1, j = 2$; степени соответствующих полиномов суть 9 и 17. Построение этих полиномов требует соответственно выполнения 220 и 1140 условий. Хотя автор и приводит связанные с этими полиномами оценки аппроксимации, однако, как он замечает сам эти полиномы практически не применимы.

3°. В статье Бергера и Маркса [1] рассмотрена аппроксимация функций класса $C^{(2)}(\Omega)$ кусочно полиномиальными функциями 5-й степени на сетке из конгруэнтных прямоугольных треугольников. Область Ω предполагается многоугольной; в каждом из треугольников сетки интерполирующий полином 5-й степени определяется значениями интерполируемой функции и ее первых производных в вершинах треугольника. Погрешность аппроксимации в $C(\Omega)$ и $C^{(1)}(\Omega)$ равна соответственно $O(h^2)$ и $O(h)$.

4°. В работе Чистова [1] рассматривается задача аппроксимации по МКЭ функции $x \in W_p^{(3)}(\Omega)$ (Ω - конечная область в R_2) на сетке, образованной непересекающимися правильными треугольниками с длиной стороны h . На такие сетки распространено определение исходных функций, а так же способ построения координатных функций из исходных; носители координатных функций суть правильные шестиугольники с центрами в вершинах треугольников сетки. Выведены фундаментальные соотношения, аналогичные соотношениям гл. III и дающие необходимые и достаточные условия полноты системы координатных функций в $W_p^{(3)}(\Omega)$. Для $j = 2$ построена система кусочно полиномиальных исходных функций 5-й степени. Эти функции зависят от существенно меньшего числа параметров, чем упомяну-

ные в п. 1⁰ полиномы Женишека, поэтому они приводят к ϵ габран-
чекским системам МКЭ более низкого порядка.

§ 6. Теорема об устойчивости.

1⁰. Устойчивость МКЭ на симплицальных сетках исследована в
работе Чистова [2] для первой краевой задачи. Рассмотрим фор-
мально самосопряженное эллиптическое уравнение

$$\sum_{|\alpha|, |\beta|=0}^s (-1)^{|\alpha|} D^{\alpha} (A_{\alpha\beta} D^{\beta} x) = f(t), \quad f \in L_2(\Omega) \quad (Y.6.1)$$

при краевом условии

$$D^{\gamma} x|_{\partial\Omega} = 0, \quad 0 \leq |\gamma| \leq s-1. \quad (Y.6.2)$$

Будем считать, что оператор задачи (Y. 6.1) - (Y.6.2) положи-
тельно определенный в $L_2(\Omega)$ и что $\Omega = \bar{\Omega}$ - конечная область в R_m .

Пусть Ω^h - многогранник, $\Omega^h \subset \bar{\Omega}$, и пусть в Ω^h выбрано
конечное множество t_k узлов, занумерованных в произвольном по-
рядке, $k = 1, 2, \dots, L$. Разобьем Ω^h на симплексы S_k , $\tau = 1, 2, \dots, \tau_0$
с вершинами в узлах t_k . Потребуем, чтобы S_k удовлет-
воряли "условию совместности" (см. Оден [1], п. 7.3). Обозна-
чим через $l_{\tau k}$ длины ребер симплекса S_k , сходящихся к вершине
с номером k . Пусть $P_{\tau k}$ - параллелепипед, построенный на этих
ребрах, и $P_{\tau k}$ - параллелепипед с ребрами единичной длины, постро-
енный на тех же направлениях.

Если Q множество в R_m , и $|Q|$ - объем (мера) этого множе-
ства, то, как известно,

$$|S_k| = \frac{|P_{\tau k}|}{m!} = \frac{|P_{\tau k}|}{m!} \prod_{i=1}^m l_{\tau k i}. \quad (Y.6.3)$$

Введем обозначения: $h_1 = \max h_{\tau k i}$, $h_2 = \min h_{\tau k i}$, $\vartheta = \min |P_{\tau k}|$.
Допустим, что

$$h_1/h_2 \leq C_1, \quad \vartheta \geq C_2; \quad C_1, C_2 = \text{const}. \quad (Y.6.4)$$

С каждым узлом t_k свяжем координатные функции $\Phi_{ik}(t)$;
здесь i - мультииндекс размерности m , пробегаящий множество
 $0 \leq |i| \leq \bar{s}$, причем число \bar{s} одинаково для всех k . Будем счи-
тать, что функции Φ_{ik} обладают следующими свойствами:

а) $\Phi_{ik} \in C^{(s-1)}(R_m) \cap W_2^{(s)}(R_m)$;

б) $\text{supp } \Phi_{ik} \subset Z_k$, где Z_k - многогранник, образова-

ный симплексами с общей вершиной t_k ;

$$в) \Phi_{ik}^{(h)}(t_k) = \delta_{di}, \quad 0 \leq |d|, |i| \leq \tilde{\delta};$$

г) Существует такая последовательность $h_n \rightarrow 0$, что последовательность соответствующих конечных сумм координатных функций $\{\Phi_{ik}\}$ полна в энергетическом пространстве H_A нашей задачи.

д) Пусть S_0 - симплекс с вершинами $(0, 0, \dots, 0)$, $(1, 0, \dots, 0)$, \dots , $(0, \dots, 0, 1)$. Существует такое аффинное преобразование симплекса S_0 на S_v : $t = \omega_v(\tau)$, что если $y_{ikt}(\tau) = \Phi_{ik}(\omega_v(\tau))$, то наименьшее собственное число матрицы скалярных произведений $\langle y_{ikt}, y_{i'kt'} \rangle_{L_2(S_0)}$ ограничено снизу положительным числом λ_0 , которое может зависеть только от области Ω и от постоянных C_1, \dots, C_n из неравенств (У.6.4).

Из условия д) вытекает, что процесс МКЭ сойдется при $h_n \rightarrow 0$ в норме H_A . Чем более он сойдется в более слабой норме $L_2(\Omega)$, поэтому, если x_* есть точное решение, а

$$x_h = \sum_{|i|=0}^{\tilde{\delta}} \sum_{k=1}^L a_{ik}^{(h)} \Phi_{ik}(t)$$

есть приближенное решение задачи (У.6.1-2), то существует такая постоянная C_0 , что

$$\|x_h^{(h)}\|_{2,\Omega} \leq C_0 \|x_h\|_{2,\Omega}. \quad (У.6.5)$$

Рассмотрим подпространство $H_A^{(h)}$ энергетического пространства H_A , образованное функциями вида

$$x(t) = \sum_{|i|=0}^{\tilde{\delta}} \sum_{k=1}^L a_{ik} \Phi_{ik}(t). \quad (У.6.6)$$

Оценим норму вектора a с составляющими a_{ik} через норму $\|x^{(h)}\|$:

$$\begin{aligned} \|x^{(h)}\|_{2,L^2}^2 &= \sum_{v=1}^q \int_{S_v} \left[\sum_{|i|=0}^{\tilde{\delta}} \sum_{k=1}^L a_{ik} \Phi_{ik}(t) \right]^2 dt = \\ &= \int_{S_0} \sum_{v=1}^q \left[\sum_{|i|=0}^{\tilde{\delta}} \sum_{k=1}^L a_{ik} y_{ikt}(\tau) \right]^2 J_v d\tau. \end{aligned} \quad (У.6.7)$$

Здесь J_v - якобиан преобразования симплекса S_v в симплекс S_0 ; как хорошо известно $|J_v| = \varrho^m |S_v|$. Далее, из соотношений (У.6.3-4) следует

$$\|x^{(h)}\|_{2, \Omega^h} \geq \frac{2^m \sigma}{m!} h_2^m (M_{h_2}^{d, d})_{R_N} \geq \frac{2^m \sigma \Lambda_0}{m!} h_2^{m\nu} \|d\|_{\dots}^2; \quad (У.6.8)$$

здесь N - число составляющих вектора d . Так как $\|x^{(h)}\|_{2, \Omega^h} \leq \|x^{(h)}\|_{2, \Omega}$, то отсюда и из (У.6.5) следует, что величина $h_2^m \|a^{(h)}\|_{R_N}^2$ ограничена. Обозначим теперь через \bar{X}_N и \bar{Y}_N гильбертовы N -мерные пространства с нормами

$$\|\cdot\|_{\bar{X}_N} = h^{m\nu/2} \|\cdot\|_{R_N}, \quad \|\cdot\|_{\bar{Y}_N} = h^{-m\nu/2} \|\cdot\|_{R_N},$$

где h любое число из последовательности $\{h_n\}$.

На основании замечания III.3.4 можно утверждать, что процесс вычисления коэффициентов $a_{ik}^{(h_n)}$ устойчив относительно погрешностей искажения в последовательности (\bar{X}_N, \bar{Y}_N) , а процесс вычисления приближенного решения $x^{(h)}$ устойчив в последовательности (\bar{X}_N, \bar{Y}_N) .

Аналогичные теоремы можно доказать при произвольных краевых условиях, оставив тех оператор (У.6.1) положительно определенным, если обычный функционал энергии заменить функционалом, который ввел Руховец [1].

ГЛАВА VI

Погрешности приближенного решения интегральных уравнений.

Проблема, формулированная в 1-м звании данной главы, по-видимому, пока еще исследована не вполне достаточно. Перечислим те результаты, которые кажутся нам более важными.

Канторович (см. Канторович и Крылов [1], гл. II) исследовал погрешность аппроксимации метода механических квадратур и метода замены ядра вырожденным для уравнений Фредгольма. Большое количество работ, посвященных приближенному решению Фредгольмовых уравнений, а также одномерных сингулярных уравнений, принадлежит Габдулхаеву (подробная библиография по ведена в его книге [3]). Некоторые результаты по приближенному решению интегральных уравнений Фредгольмовых и сингулярных одномерных, приведены в книге автора и Смолицкого [1], глава VI. В книге Гохберга и Фельдмана [1] с большой полнотой изложено применение проекционных методов к "уравнениям в свертках", т.е., к уравнениям, являющимся обобщением известных уравнений Винера - Хоупа. Некоторые новые подходы к приближенному решению уравнений Фредгольма указаны в книге Анселона [1]. Приближенному решению одномерных сингулярных интегральных уравнений посвящена книга Пресдорфа и Зильсэрманн [1]. Ряд результатов по приближенному решению сингулярных интегральных уравнений, одномерных и многомерных, содержится в книге автора - Пресдорфа [1]. Приближенному решению многомерных сингулярных интегральных уравнений посвящены работы автора и Радевой [1], автора [25]. Отметим еще сборники статей под редакцией Делвса и Уолша [1] - Андерссена, де-Хугта Липаса [1], посвященные приближенному решению интегральных уравнений, а так же весьма интересную статью Пресдорфа и Шмидта [2], в которой получено необходимое и достаточное условие сходимости одного варианта метода коллокации, связанного с МКЭ, и дана оценка соответствующей погрешности аппроксимации. В работах автора [18-21, 23] развит метод приближенного вычисления резольвенты Фредгольма, основанный на конформной аппроксимации ядра.

Для упрощения записи, а отчасти и выкладок, мы будем рассматривать ниже скалярные уравнения в R_1 и будем считать, что промежуток контура решения есть отрезок $[0, 1]$. Распространение результатов на более общие уравнения Фредгольма очевидно. Ими

чином, пожалуй, является аппроксимация ядра по М(Э в том случае, когда множество интегрирования не допускает взаимно-однозначного отображения на куб; этот случай рассмотрен в статье автора [20].

§ I. Применение метода механических квадратур к уравнениям Фредгольма

Г^с. Рассмотрим уравнение

$$x(t) - \int_0^1 K(t, \tau) x(\tau) d\tau = f(t). \quad (VI.1.1)$$

Допустим, что $f \in C^{(s)}[0, 1]$, $K \in C^{(s, s)}([0, 1] \times [0, 1])$ и что число 1 не характеристическое для ядра $K(t, \tau)$. Условие на ядро надо понимать так: в квадрате $[0, 1] \times [0, 1]$ существуют и непрерывны все производные вида

$$\frac{\partial^{k+l} K}{\partial t^k \partial \tau^l}; \quad 0 \leq k, l \leq s.$$

Из наших условий вытекает, что уравнение (VI.1.1) имеет одно и только одно решение $x_s \in C^{(s)}[0, 1]$.

Будем приближенно решать уравнение (VI.1.1) по методу механических квадратур, заменяя интеграл по квадратурной формуле вида

$$\int_0^1 y(\tau) d\tau = \sum_{k=1}^n C_k^{(n)} y(\tau_k^{(n)}) + f_n. \quad (VI.1.2)$$

Примем, что эта формула обладает следующими свойствами:

$$C_k^{(n)} \geq 0, \quad \sum_{k=1}^n C_k^{(n)} = 1; \quad (VI.1.5)$$

если $y \in C^{(s)}[0, 1]$, то

$$\rho_n \leq C_3 h^s \|y^{(s)}\|_C; \quad h = \max_k (\tau_{k+1}^{(n)} - \tau_k^{(n)}); \quad C_3 = \text{const}. \quad (VI.1.4)$$

Заменяя в (VI.1.1) интеграл по формуле (VI.1.2) и отбрасывая остаточный член, получим новое уравнение

$$y^{(n)}(t) = \sum_{k=1}^n C_k^{(n)} K(t, \tau_k^{(n)}) y^{(n)}(\tau_k^{(n)}) + f(t), \quad \tau_k = \tau_k^{(n)} \quad (VI.1.5)$$

Решая $y_j^{(n)}(t)$ которого будем рассматривать как приближенное решение уравнения (VI.1.1). Чтобы определить $y_j^{(n)}(\tau_j^{(n)})$, положим в (VI.1.1) $t = \tau_j$, $1 \leq j \leq n$. Это даст нам алгебраическую систему

$$y_j^{(n)} - \sum_{k=1}^n c_k^{(n)} \mathcal{K}(\tau_j, \tau_k) y_k^{(n)} = f_j^{(n)}, \quad 1 \leq j \leq n, \quad (VI.1.6)$$

где $y_j^{(n)} = y^{(n)}(\tau_j)$, $f_j^{(n)} = f(\tau_j)$. Если система (VI.1.5) разрешима единственным образом, то тем же свойством обладает уравнение (VI.1.5). Очевидно и обратное; очевидно также, что любое решение уравнения (VI.1.5) принадлежит к классу $C^{(3)}[0, 1]$.

Операторы

$$\int_0^1 \mathcal{K}(t, \tau) x(\tau) d\tau, \quad \sum_{k=1}^n c_k^{(n)} \mathcal{K}(t, \tau_k) x(\tau_k) \quad (VI.1.7)$$

будем рассматривать как операторы в $C^{(3)}[0, 1]$; оценим норму разности. Эта норма равна

$$\sup_{\|x\|=1} \left\| \int_0^1 \frac{\partial^j \mathcal{K}(t, \tau)}{\partial t^j} x(\tau) d\tau - \sum_{k=1}^n c_k^{(n)} \frac{\partial^j \mathcal{K}(t, \tau_k)}{\partial t^j} x(\tau_k) \right\|_C.$$

По оценке (VI.1.4) это не превосходит величины

$$\sup_{\|x\|=1} C h^3 \sum_{j=0}^2 \left\| \frac{\partial^3}{\partial t^3} \frac{\partial^j \mathcal{K}(t, \tau)}{\partial t^j} x(\tau) d\tau \right\|_C \leq C h^3; \quad C = \text{const}.$$

Таким образом, при достаточно малом h разность операторов (VI.1.7) сколь угодно мала по норме $C^{(3)}$.

Интегральные операторы в (VI.1.7) обозначим соответственно через \mathcal{K} и \mathcal{K}_n , а их разность через \mathcal{G}_n . Уравнение (VI.1.5) можно записать в виде

$$y^{(n)}(t) - (\mathcal{K} y^{(n)})(t) + (\mathcal{G}_n y^{(n)})(t) = f(t).$$

Отсюда

$$y^{(n)}(t) + ((I - \mathcal{K})^{-1} \mathcal{G}_n y^{(n)})(t) = (I - \mathcal{K})^{-1} f(t) = \alpha_*(t)$$

и, следовательно,

$$\|\alpha_* - y_*^{(n)}\|_{C^{(3)}} \leq \|(I - \mathcal{K})^{-1}\|_{C^{(3)}} \cdot \|\mathcal{G}_n\|_{C^{(3)}} \cdot \|y_*^{(n)}\|_{C^{(3)}}.$$

При достаточно малом h будет $\|(I - \mathcal{K})^{-1}\|_{C^{(3)}} \cdot \|\mathcal{G}_n\|_{C^{(3)}} < 1/2$

$$\|y_*^{(n)}\|_{C^{(3)}} \cdot \|\alpha_*\|_{C^{(3)}} \leq \frac{1}{2} \|y_*^{(n)}\|_{C^{(3)}}; \quad \|y_*^{(n)}\|_{C^{(3)}} \leq 2 \|\alpha_*\|_{C^{(3)}}.$$

Теперь,

$$\|\alpha_* - y_*^{(n)}\|_{C^{(3)}} \leq C h^3 \|\alpha_*\|_{C^{(3)}}, \quad C = \text{const}.$$

(VI.1.8)

остоянную C нетрудно уточнить:

$$C = 2C_3 \| (I - \mathcal{K})^{-1} \|_{C(\Omega)}. \quad (VI.1.9)$$

Формулы (VI.1.8) и (VI.1.9) дают оценку по точности аппроксимации метода механических квадратур для уравнения Фредгольма.

2°. Рассмотрим случай, когда ядро и свободный член уравнения (VI.1.1) только непрерывны. Для произвольной функции $u(t)$ определим ее модуль непрерывности $\omega(u, h)$ и допустим, что $|u(t') - u(t'')| \leq M_u r(h)$, где $|t' - t''| \leq h$, $r(h) \xrightarrow{h \rightarrow 0} 0$, причем M_u зависит от u , но не от h . Наименьшее возможное значение M_u обозначим через $M(u)$:

$$\omega(u, h) \leq M(u)r(h). \quad (VI.1.10)$$

Допустим, что модули непрерывности ядра $\mathcal{K}(t, \tau)$ и свободного члена $f(t)$ удовлетворяют неравенству (VI.1.10). Тогда тому же неравенству удовлетворяет и решение уравнения (VI.1.1).

Введем пространство S функций, непрерывных на отрезке $[0, 1]$ с модулем непрерывности, удовлетворяющим неравенству (VI.1.11) и с нормой

$$\| \cdot \|_S = \| \cdot \|_C + M(\cdot). \quad (VI.1.11)$$

Очевидно, $f, x_n \in S$. Докажем, что операторы (VI.1.7) ограничены в S . Обозначим $(\mathcal{K}x)(t) = y(t)$. Имеем $\|y\|_C \leq \|\mathcal{K}\|_C \cdot \|x\|_C$. Далее,

$$\omega(y, h) \leq \omega(\mathcal{K}, h) \|x\|_C \leq M(\mathcal{K})r(h) \|x\|_C. \quad (VI.1.12)$$

Отсюда $M(y) \leq M(\mathcal{K}) \|x\|_C$ и

$$\|y\|_S \leq [\|\mathcal{K}\|_C + M(\mathcal{K})] \cdot \|x\|_C \leq [\|\mathcal{K}\|_C + M(\mathcal{K})] \cdot \|x\|_S. \quad (VI.1.13)$$

Таким образом, $\|\mathcal{K}\|_{S \rightarrow S} \leq \|\mathcal{K}\|_C + M(\mathcal{K})$. Теперь рассмотрим оператор \mathcal{K}_n . Обозначим $\psi_n(t) = (\mathcal{K}_n x)(t)$. Тогда

$$|\psi_n(t)| \leq \sum_{k=1}^n c_k^{(n)} \|\mathcal{K}\|_C \cdot \|x\|_C = \|\mathcal{K}\|_C \cdot \|x\|_C;$$

$$\omega(\psi_n, h) \leq \sum_{k=1}^n c_k^{(n)} \omega(\mathcal{K}, h) \|x\|_C = \omega(\mathcal{K}, h) \|x\|_C \leq M(\mathcal{K}) \|x\|_C r(h).$$

Отсюда $M(\psi_n) \leq M(\mathcal{K}) \|x\|_C$ и $\|\mathcal{K}_n\|_{S \rightarrow S} \leq \|\mathcal{K}\|_C + M(\mathcal{K})$.

Допустим, что для функций $y(t)$, удовлетворяющих условиям

(VI.1.10), остаточный член квадратурной формулы (VI.1.4) имеет оценку

$$p_n \leq C_0 r(h) [\|y\|_C + M y]; \quad C_0 = \text{const}, \quad r(h) \xrightarrow{h \rightarrow 0} 0. \quad (\text{VI.1.14})$$

Дополнительно предположим, что

$$M_y(\mathcal{K}) := \frac{1}{r(h)} \sup_{|t'-t''| \leq h} |\mathcal{K}(t'', \tau) - \mathcal{K}(t', \tau)| \leq C_1;$$

$$M_{t,\tau}(\mathcal{K}) := \frac{1}{r^2(h)} \sup_{\substack{|t'-t''| \leq h \\ |\tau'-\tau''| \leq h}} |\mathcal{K}(t'', \tau'') - \mathcal{K}(t'', \tau') - \\ - \mathcal{K}(t', \tau'') + \mathcal{K}(t', \tau')| \leq C_1; \quad C_1 = \text{const}.$$

Оценим $\|\sigma_n\|_{S \rightarrow S}$. Имеем $\|\sigma_n x\|_S = \|\sigma_n x\|_C + M(\sigma_n x)$.
 По оценке (VI.1.14)

$$\|\sigma_n x\|_C = \|\mathcal{K}x - \mathcal{K}_n x\|_C \leq C_0 r(h) [\|\mathcal{K}\|_C \cdot \|x\|_C + M(\mathcal{K}x)] \leq \\ \leq 2C_0 \|\mathcal{K}\|_C \cdot \|x\|_C r(h) \leq C_2 \|x\|_C r(h); \quad C_2 = \text{const}. \quad (\text{VI.1.15})$$

Далее, по той же оценке (VI.1.14)

$$\omega(\sigma_n x, h) = \sup_{|t'-t''| \leq h} \left| \int_{t'}^{t''} [\mathcal{K}(t'', \tau) - \mathcal{K}(t', \tau)] x(\tau) d\tau - \right.$$

$$\left. - \sum_{k=1}^n C_k^{(n)} [\mathcal{K}(t'', \tau_k) - \mathcal{K}(t', \tau_k)] x(\tau_k) \right| \leq$$

$$\leq C_0 r(h) [\max_{\tau} |\mathcal{K}(t'', \tau) - \mathcal{K}(t', \tau)| \cdot \|x\|_C +$$

$$+ M_{t,\tau}(\mathcal{K}(t'', \tau) - \mathcal{K}(t', \tau)) \|x\|_C] \leq 2C_0 C_1 r^2(h) \|x\|_C.$$

Отсюда $M(\sigma_n x) \leq 2C_0 C_1 r(h) \|\omega\|_C$ и $\|\sigma_n x\|_S \leq Cr(h) \|x\|_C \leq Cr(h) \|x\|_S$, $C = \text{const}$. Таким образом, $\|\sigma_n\|_{S \rightarrow S} \leq Cr(h)$ и, при h достаточно малом, будет $\|(I - \mathcal{K})\|_{S \rightarrow S} \|\sigma_n\|_{S \rightarrow S} < \frac{1}{2}$.

Поступая далее, как в конце п. I^o, получим

$$\|x_n - y_n^{(n)}\| \leq Cr(h), \quad C = \text{const}. \quad (\text{VI.1.16})$$

3^o. Рассмотрим погрешность искажения для уравнений Фредгольма. Для удобства остановимся на случае п. I^o. Обозначим

через A_n матрицу системы (VI.1.6): $A_n = I_n - \mathcal{K}^{(n)}$, где $\mathcal{K}^{(n)}$ - матрица элементов $\mathcal{K}_k^{(n)}(\tau_j, \tau_k)$, I_n - единичная матрица порядка $n \times n$. Будем рассматривать $\vec{y}^{(n)} = (y_1^{(n)}, y_2^{(n)}, \dots, y_n^{(n)})$ и $\vec{f}^{(n)} = (f_1, f_2, \dots, f_n)$ как векторы в пространстве U_n (см. § 9 гл. III). Пусть $\delta_n^{(0)}$ - наименьшее сингулярное число матрицы A_n . Так как $y_k^{(n)}$ суть значения y_k^* - решения уравнения (VI.1.5) - в точках τ_k , а $y_k^* \xrightarrow{C} x_k$ и тем более $y_k^* \xrightarrow{L_2} x_k$, то нормы $\|\vec{y}^{(n)}\|_{U_n}$ ограничены в совокупности и по следствию III.2.2 для устойчивости процесса (VI.1.6) в последовательности (U_n, U_{n+1}) необходимо и достаточно, чтобы $\|A_n\|_{U_n \rightarrow U_n} \rightarrow y_n \leq C = \text{const}$.

Но

$$\|A_n^{-1}\|_{U_n \rightarrow U_n} = \|A_n^{-1}\|_{R_n \rightarrow R_n}$$

и для устойчивости процесса (VI.1.6) в упомянутом явном смысле необходимо и достаточно, чтобы $\delta_n^{(0)}$ было положительно ограничено снизу. Легко убедиться, что это же условие необходимо и достаточно для устойчивости процесса вычисления приближенного решения в соответствующей последовательности пар подпространств.

4°. В общем случае можно оценить погрешности искажения и округления для метода механических квадратур, используя соответствующие оценки для линейных алгебраических систем (см. § 4 гл. I).

§ 2. Уравнения Эрдганья, разномне и зрашлым

1°. Пусть в уравнении (VI.1.7) ядро удовлетворяет неравенству

$$\max_{0 < t, \tau < 1} |\mathcal{K}(t, \tau)| = q < 1. \quad (\text{VI.2.1})$$

Тогда решение этого уравнения можно построить по методу простой итерации:

$$x_n = f(t) + \sum_{k=1}^{\infty} \int_0^1 \mathcal{K}_n(t, \tau) f(\tau) d\tau; \quad (\text{VI.2.2})$$

\mathcal{K}_n - втерироганные ядра. Ряд (VI.2.2) сходится в $C[0, 1]$ как прогрессия со знаменателем q , точнее, n -й член этого ряда не превосходит величины $q^n \int_0^1 |f(\tau) d\tau|$. Удобнее проводить вычисления не по формуле (VI.2.2), а по обычной схеме итераций:

$$x_0(t) = f(t); \quad n \geq 1, \quad x_n(t) = f(t) + \int_0^1 \mathcal{K}(t, \tau) x_{n-1}(\tau) d\tau. \quad (\text{VI.2.3})$$

Однако, если ядро $\mathcal{K}(t, \tau)$ имеет более или менее сложную структуру, то вычисление интегралов, входящих в формулу (VI.2.3), становится затруднительным. Поэтому мы прибегнем к следующему приему (см. работы автора [18, 22], где этот прием использован в более общей ситуации).

Введем функцию

$$\omega(t) = \begin{cases} t, & 0 \leq t \leq 1, \\ 2-t, & 1 < t \leq 2, \\ 0, & t \notin [0, 2] \end{cases} \quad (\text{VI.2.4})$$

и обозначим

$$\mathcal{K}^h(t, \nu) = \sum_{j, k=1}^{2n-1} \mathcal{K}(t_{j_n}, t_{k_n}) \omega\left(\frac{t}{h} - j\right) \omega\left(\frac{\nu}{h} - k\right). \quad (\text{VI.2.5})$$

Здесь n — произвольно выбранное натуральное число, $h = 1/2n$, $t_j = (j+1)h$. Если $\mathcal{K} \in C^2(Q)$, $Q = [0, 1] \times [0, 1]$, то (см. работу автора [22], гл. III, § I)

$$\|\mathcal{K} - \mathcal{K}^h\|_{C(Q)} \leq C \max_{|\alpha|=2} \|\mathcal{D}^\alpha \mathcal{K}\|_{C(Q)} h^2.$$

Постоянную C в последнем неравенстве можно вычислить по формуле (9.14) статьи автора [24]; см. также книгу автора [3] Кар. IV, § 9). Заметим, что в нашем случае следует в упомянутой формуле (9.14) положить $\delta = 1$, а тогда легко получается $C = 1$ так что

$$\|\mathcal{K} - \mathcal{K}^h\|_{C(Q)} \leq 10 \max_{|\alpha|=2} \|\mathcal{D}^\alpha \mathcal{K}\|_{C(Q)} h^2. \quad (\text{VI.2.6})$$

Возьмем n столь большим, чтобы правая часть оценки (VI.2.6) была меньше, чем $1 - q$. Тогда, очевидно, $\|\mathcal{K}^h\|_{C(Q)} := q_n < 1$. Заменяя в уравнении (VI.1.1) ядро \mathcal{K} на \mathcal{K}^h , а свободный член $f(t)$ — его конечноэлементной аппроксимацией

$$f^h(t) = \sum_{j=1}^{2n-1} f(t_{j_n}) \omega\left(\frac{t}{h} - j\right).$$

Равность $f(t) - f^h(t)$ есть разность ординат кривой $y = f(t)$ и вписанной в нее ломаной $y = f^h(t)$. Если $f \in C^{(2)}[0, 1]$, то эту разность можно вычислить по формуле

$$y_{j+1} - y_{j+2}, \quad f(t) - f^h(t) = -\frac{1}{h} \int_{t_{j_n}}^{t_{j+1}_n} G(t - t_{j_n}, \tau - t_{j+1}_n) f''(\tau) d\tau,$$

где

$$G(t, \tau) = \begin{cases} t(h-\tau), & 0 \leq t \leq \tau, \\ \tau(h-t), & 0 \leq \tau \leq t. \end{cases}$$

отсюда

$$\|f - f^{(h)}\|_C \leq h^2 \|f''\|_C. \quad (\text{VI.2.7})$$

2°. Оценим погрешность аппроксимации от такой замены. Новое интегральное уравнение имеет вид

$$x^{(h)}(t) - \int_0^1 \mathcal{K}^h(t, \tau) x^{(h)}(\tau) d\tau = f^{(h)}(t). \quad (\text{VI.2.8})$$

Сравнивая это с общей схемой п.4° § 4 гл. I, видим, что в данном случае $\Gamma = \mathcal{K} - \mathcal{K}^h$, где \mathcal{K} и \mathcal{K}^h - интегральные операторы с ядрами $\mathcal{K}(t, \tau)$ и $\mathcal{K}^h(t, \tau)$ соответственно.

По формуле (I.4.14)

$$\|x - x^{(h)}\| \leq \frac{\|(I - \mathcal{K})^{-1}\|}{1 - \|(I - \mathcal{K})^{-1}\| \|\mathcal{K} - \mathcal{K}^h\|} [\|\mathcal{K} - \mathcal{K}^h\| \cdot \|x\| + \|\delta\|];$$

все нормы взяты в $C[0, 1]$ или, соответственно, в $C([0, 1] \times [0, 1])$.

Используя неравенства (VI.2.6) и (VI.2.7), а также неравенство

$\|(I - \mathcal{K})^{-1}\| \leq (1 - q)^{-1}$, получим оценку погрешности аппроксимации:

$$\begin{aligned} \|x - x^{(h)}\| &\leq \\ &\leq [1 - 10q \max_{|k|=2} \|\mathcal{D}^k \mathcal{K}\| h^2]^{-1} [10 \max_{|k|=2} \|\mathcal{D}^k \mathcal{K}\| \|x\| + \|f''\|] h^2; \quad (\text{VI.2.9}) \end{aligned}$$

можно еще воспользоваться тем, что

$$\|x\| \leq \frac{\|f\|}{1 - q}. \quad (\text{VI.2.10})$$

3°. Решим уравнение (VI.2.8) по методу итераций и оценим погрешность алгоритма для этого уравнения. В данном случае итерации сходятся как прогрессия с знаменателем q_1 . Если выполнять точно N итераций, то погрешность алгоритма будет иметь оценку

$$\|x^{(h)} - x_N^{(h)}\| \leq$$

$$\leq \sum_{k=N+1}^{\infty} \left\| \int_0^1 \mathcal{K}_k^h(t, \tau) f^{(h)}(\tau) d\tau \right\| \leq [\|f\| + h^2 \|f''\|] \frac{q_1^{N+1}}{1 - q_1}. \quad (\text{VI.2.11})$$

4°. Допустим, что числа $\mathcal{K}(t_j, t_k)$ и $f(t_j)$ вычислены с погрешностями, тогда мы имеем точную формулу

$$x_k^{(h)} = f^{(h)}(t_k) + \int_0^1 \mathcal{K}^h(t, \tau) x_{k-1}^{(h)}(\tau) d\tau. \quad (VI.2.12)$$

Нетрудно видеть, что

$$x_{k-1}^{(h)}(t) = \sum_{j=0}^{2n-1} a_{j,k-1}^{(h)} \omega\left(\frac{t}{h} - j\right).$$

Подставим это, а также значения $f^{(h)}(t)$ и $\mathcal{K}^h(t, \tau)$, в (VI.2.12) и приравняв коэффициенты при $\omega\left(\frac{t}{h} - j\right)$, получим

$$a_{j,k}^{(h)} = f(t_{j+1}) + \sum_{l,m=0}^{2n-1} a_{m,k-1} \mathcal{K}(t_{j+1}, t_{l+1}) \int_0^1 \omega\left(\frac{\tau}{h} - l\right) \omega\left(\frac{\tau}{h} - m\right) d\tau.$$

Интеграл в последнем тождестве просто вычисляется: обозначив через $h\omega_{lm}$, имеем $\omega_{lm} = \omega_{ml}$ и

$$\omega_{lm} = \begin{cases} 0, & |m-l| \geq 2 \\ 1/6, & m-l = \pm 1 \\ 2/3, & m=l, 0 \leq m \leq 2n-2 \\ 1/3, & m=l=2n-1. \end{cases} \quad (VI.2.13)$$

В результате мы приходим к системе рекуррентных соотношений коэффициентов $a_{jk}^{(h)}$:

$$a_{jk}^{(h)} = f(t_{j+1}) + h \sum_{l,m=0}^{2n-1} a_{m,k-1} \mathcal{K}(t_{j+1}, t_{l+1}) \omega_{lm}; \quad (VI.2.14)$$

систему (VI.2.14) надо решать при начальном условии $a_{j_0}^{(h)} = f(t_{j_0})$.

Если мы введем в пространство $2n$ -компонентных векторов \mathcal{A} по формуле $\|\mathcal{A}\| = \max |a_i|$, то норма оператора в правой части равенства (VI.2.14) равна

$$B = h \max_{l,m=0}^{2n-1} |\mathcal{K}(t_{j+1}, t_{l+1})| \omega_{lm}. \quad (VI.2.15)$$

Если в дальнейшем значения чисел ω_{lm} (формула (VI.2.13)), считаем, что

$$\sum_{l,m=0}^{2n-1} \omega_{lm} |\mathcal{K}(t_{j+1}, t_{l+1})| \approx \frac{1}{h} \max |\mathcal{K}(t, \tau)| = q \cdot 2n.$$

Отсюда $B \leq h \cdot 2nq = q$, и при фиксированном индексе j итерационный процесс (VI.2.14) сходится. Если ограничиться N итерациями, то погрешность алгоритма для процесса (VI.2.14) оценивается величиной

$$\frac{q^{N+1}}{1-q} \|f(t)\|_C. \tag{VI.2.16}$$

5°. В нашей задаче можно не выделять погрешность искажения и считать, что остальные погрешности вызваны округлением. Погрешность округления легко оценить по способу § 3 гл. IV; так как итерационный процесс (VI.2.14) сходится, как геометрическая прогрессия, то при бесконечном возрастании N погрешность округления остается ограниченной.

§ 3. Погрешность резольвенты Фредгольма

1°. Результаты настоящего параграфа содержатся в статье автора [18], ряд обобщений - в его статьях [19-21]; материал этих статей изложен в книге автора [22]. Ввиду этого мы здесь ограничимся только постановкой задачи и формулировкой результатов, и притом для простейшего класса интегральных уравнений и простейшей аппроксимации ядра. Подробное доказательство мы проведем только для оценки погрешности алгоритма, которая уточняет соответствующую оценку из [16] и [22].

2°. Рассмотрим интегральное уравнение Фредгольма

$$x(t) - \lambda \int_0^1 K(t, \tau) x(\tau) d\tau = f(t) \tag{VI.3.1}$$

с ядром $K \in C^{(2)}$; λ - произвольный, вообще говоря, комплексный параметр. Требование на ядро можно ослабить - достаточно считать его непрерывным, но при этом некоторые последующие оценки окажутся хуже. Будем считать еще, что $|K(t, \tau)| < 1$ - этого всегда можно добиться подходящим изменением параметра.

Напомним основные факты, относящиеся к резольвенте Фредгольма.

1. Резольвента Фредгольма $\Gamma(t, \tau; \lambda)$ для ядра $K(t, \tau)$ имеет вид

$$\Gamma(t, \tau; \lambda) = \frac{D(t, \tau; \lambda)}{D(\lambda)}, \tag{VI.3.2}$$

где "определитель Фредгольма" $D(\lambda)$ и "первый минор Фредгольма" $D(t, \tau; \lambda)$ - определенные функции от λ .

2. Если $D(\lambda) \neq 0$, то уравнение (VI.3.1) имеет один и тот же...

но от решения

$$x(t) = f(t) + \lambda \int_0^1 \Gamma(t, \tau; \lambda) f(\tau) d\tau. \quad (VI.3)$$

3. Коэффициенты рядов

$$D(\lambda) = \sum_{k=0}^{\infty} \frac{(-1)^k}{k!} c_k \lambda^k \quad (VI.3)$$

и

$$D(t, \tau; \lambda) = \sum_{k=0}^{\infty} \frac{(-1)^k}{k!} B_k(t, \tau) \lambda^k \quad (VI.3)$$

удовлетворяют рекуррентным соотношениям

$$c_0 = 1, \quad c_k = \int_0^1 B_{k-1}(\tau, \tau) d\tau; \quad (VI.3)$$

$$B_k(t, \tau) = \mathcal{K}(t, \tau),$$

$$B_k(\tau, \tau) = c_k \mathcal{K}(\tau, \tau) - k \int_0^1 \mathcal{K}(\tau, \tau') B_{k-1}(\tau', \tau) d\tau'. \quad (VI.3)$$

4. Справедлива формула, выражающая коэффициенты $B_k(t, \tau)$ непосредственно через $\mathcal{K}(t, \tau)$:

$$k > 0, \quad B_k(t, \tau) = \int_0^1 \dots \int_0^1 \begin{vmatrix} \mathcal{K}(t, \tau) & \mathcal{K}(t, \tau_1) & \dots & \mathcal{K}(t, \tau_k) \\ \mathcal{K}(\tau_1, \tau) & \mathcal{K}(\tau_1, \tau_1) & \dots & \mathcal{K}(\tau_1, \tau_k) \\ \dots & \dots & \dots & \dots \\ \mathcal{K}(\tau_k, \tau) & \mathcal{K}(\tau_k, \tau_1) & \dots & \mathcal{K}(\tau_k, \tau_k) \end{vmatrix} d\tau_1 \dots d\tau_k. \quad (VI.3)$$

3°. Заменим ядро $\mathcal{K}(t, \tau)$ его аппроксимацией (VI.2.5) с помощью неравенств (VI.2.6) при достаточно малом h будем

$|\mathcal{K}^h(t, \tau)| < 1$. Обозначим через $b_{j\ell}^h$ и $B_{j\ell}^h(t, \tau)$ коэффициенты рядов Фурье ядра $\mathcal{K}^h(t, \tau)$. Нетрудно убедиться, что

$$B_{j\ell}^h(t, \tau) = \sum_{k=0}^{\infty} b_{j\ell}^h \omega_{j\ell}^{(k)} \left(\frac{t}{h} - j \right) \omega_{j\ell}^{(k)} \left(\frac{\tau}{h} - \ell \right), \quad (VI.3)$$

а рекуррентные соотношения (VI.3.6), (VI.3.7) переходят в следующие:

$$c_0^h = 1, \quad c_k^h = h \sum_{j, \ell=0}^{2n-1} b_{j\ell}^h \omega_{j\ell}^{(k-1)}. \quad (VI.3)$$

$$b_{jkh}^{(k)} = \mathcal{K}((j+1)h, (l+1)h); k > 0,$$

$$b_{jkh}^{(0)} = c_k^h b_{jkh}^{(0)} - kh \sum_{l,j=1}^{2n-1} b_{jkh}^{(0)} b_{lkh}^{(k-1)} \sigma_{j+l} \quad (\text{VI.3.II})$$

Значения $W_{\nu_0}^h$ даны в формуле (VI.2.I3). Соотношения (VI.3.I0-II) позволяют вычислить резольвенту Фредгольма для ядра $\mathcal{K}^h(t, \tau)$.

4^o. В обозначениях § I гл. I \mathcal{D} есть совокупность рядов Фредгольма $\mathcal{D}(\lambda)$ и $\mathcal{D}(t, \tau; \lambda)$, \mathcal{K} - ядро $\mathcal{K}(t, \tau)$, соотношение \mathcal{U} определено равенствами (VI.3.4-7). Далее, $\mathcal{K}^{(n)}$ есть ядро $\mathcal{K}^h(t, \tau)$, $\mathcal{D}^{(n)}$ - совокупность рядов Фредгольма этого ядра, $t_n = \tau$. Погрешность аппроксимации оценивается формула:

$$|\mathcal{D}(t, \tau; \lambda) - \mathcal{D}^h(t, \tau; \lambda)| \leq ch^2 [P_3(|\lambda|^2) + |\lambda|^2 \exp(e|\lambda|^2)]^{1/2}$$

$$|\mathcal{D}(\lambda) - \mathcal{D}^h(\lambda)| \leq ch^2 [P_3(|\lambda|^2) + |\lambda|^2 \exp(e|\lambda|^2)]^{1/2}; \quad (\text{VI.3.I2})$$

P_3 - некоторый полином 3-й степ. и.

Погрешность искажения в данном случае определяется погрешностями вычисления значений $\mathcal{K}(t_{j+1}, t_{l+1})$, т.е., значений коэффициентов $b_{jkh}^{(k)}$. Поэтому нет нужды особо выделять погрешность искажения, а можно оценить оуммарно погрешности искажения и округления для процесса (VI.3.I0-II).

Оценим погрешность алгоритма. Прежде всего заметим, что ядро $\mathcal{K}^h(t, \tau)$ - вырожденное:

$$\mathcal{K}^h(t, \tau) = \sum_{j,l=1}^{2n-1} \mathcal{K}(t_{j+1}, t_{l+1}) \omega\left(\frac{t}{h} - j\right) \omega\left(\frac{\tau}{h} - l\right) = \sum_{j=1}^{2n-1} a_j(t) b_j(\tau),$$

где

$$a_j(t) = \omega\left(\frac{t}{h} - j\right), \quad b_j(\tau) = \sum_{l=1}^{2n-1} \mathcal{K}(t_{j+1}, t_{l+1}) \omega\left(\frac{\tau}{h} - l\right).$$

Отсюда следует, что ряды $\mathcal{D}^h(t, \tau; \lambda)$ и $\mathcal{D}^h(\lambda)$ суть в явном от λ степени не выше $2n$, так что $B_x^h(t, \tau) = 0$ и $C_x^h = 0, k > 2n$, и соотношениями (VI.3.I0-II) достаточно воспользоваться $2n$ раз.

Допустим теперь, что мы ограничимся вычислением коэффициентов B_x^h и C_x^h для $k \leq N, N < 2n$. Погрешность для $\mathcal{D}^h(t, \tau; \lambda)$ при этом не превосходит 1 единицы

$$\varrho_{\mathcal{N},h} = \sum_{k=\mathcal{N}+1}^{2\mathcal{N}} \frac{|B_k^h(t, \tau)|}{k!} |\lambda|^k \leq \sum_{k=\mathcal{N}+1}^{2\mathcal{N}} \frac{(k+1)^{k+1/2}}{k!} |\lambda|^k.$$

Оценим общий член последней суммы. По формуле Стирлинга

$$k! \sim \sqrt{2\pi k} \left(\frac{k}{e}\right)^k e^{\theta/2k}, \quad 0 < \theta < 1, \quad \text{Отсюда}$$

$$\frac{|\lambda|^k (k+1)^{(k+1)/2}}{k!} < \frac{\sqrt{(k+1)!} e^{k+1}}{\sqrt{2\pi(k+1)} k!} |\lambda|^k = \frac{\sqrt{e}}{\sqrt{2\pi}} \frac{(\lambda \sqrt{e})^k}{\sqrt{k!}} (\dots)^{1/4}$$

и, следовательно,

$$\varrho_{\mathcal{N},h} \leq \frac{\sqrt{e}}{\sqrt{2\pi}} \sum_{k=\mathcal{N}+1}^{2\mathcal{N}} \frac{(\lambda \sqrt{e})^k}{\sqrt{k!}} (\dots)^{1/4} < \frac{\sqrt{e}}{\sqrt{2\pi}} \sum_{k=\mathcal{N}+1}^{2\mathcal{N}} \frac{k^{-1/4} (\lambda^2 e)^{k/2}}{\sqrt{(k-1)!}} <$$

$$< \frac{\sqrt{e}}{\sqrt{2\pi}} \frac{1}{\sqrt{\mathcal{N}}} \sum_{k=\mathcal{N}+1}^{2\mathcal{N}} \frac{(\lambda^2 e)^{k/2}}{\sqrt{(k-1)!}} =$$

$$= \frac{\sqrt{e}}{\sqrt{2\pi}} \frac{1}{\sqrt{\mathcal{N}+1}} \frac{(\lambda^2 e)^{\mathcal{N}+1/2}}{\sqrt{\mathcal{N}!}} \sum_{k=\mathcal{N}+1}^{2\mathcal{N}} \frac{(\lambda^2 e)^{(k-\mathcal{N}-1)/2}}{\sqrt{(\mathcal{N}+1)(\mathcal{N}+2)\dots(k-1)}} =$$

$$\leq \frac{\sqrt{e}}{\sqrt{2\pi(\mathcal{N}+1)}} \frac{(\lambda^2 e)^{\mathcal{N}+1/2}}{\sqrt{\mathcal{N}!}} \sqrt{e^{-\mathcal{N}}} \left[\sum_{\nu=0}^{\infty} \frac{(\lambda^2 e)^{\nu}}{(\mathcal{N}+1)\dots(\mathcal{N}+\nu)} \right]^{1/2}.$$

ли $\lambda = 0$, то знаменатель в последней сумме надо заменить единицей. Заметим, что общий член этой суммы равен

$$\frac{(\lambda^2 e)^{\nu}}{\nu!} \frac{\nu!}{(\mathcal{N}+1)\dots(\mathcal{N}+\nu)} < \frac{(\lambda^2 e)^{\nu}}{\nu!} \frac{1}{\mathcal{N}+1}$$

и упомянутая сумма имеет оценку

$$1 + \sum_{\nu=1}^{\infty} \frac{(\lambda^2 e)^{\nu}}{\nu!} \frac{1}{\mathcal{N}+1} < 1 + \frac{e\lambda^2 \lambda^2 e}{\mathcal{N}+1}.$$

Окучательно, по решению алгоритма для $\mathcal{D}(t, \tau; \lambda)$ имеет оценку

$$\varrho_{\mathcal{N},h} < \frac{\sqrt{2(2\mathcal{N}-\mathcal{N})}}{\sqrt{2\pi(\mathcal{N}+1)}} \frac{(\lambda^2 e)^{\mathcal{N}+1/2}}{\sqrt{\mathcal{N}!}} \sqrt{1 + \frac{e\lambda^2 \lambda^2 e}{\mathcal{N}+1}}. \quad (\text{VI.3.13})$$

Та же величина оценивает погрешность алгоритма и для $\mathcal{D}(\lambda)$.

Формула (VI.3.13) точнее аналогичной формулы, полученной в [13, 22].

5°. Обозначим

$$\frac{\beta_{jkh}^{(k)}}{k!} = \delta_{jkh}^{(k)}; \quad \frac{c_k^h}{k!} = \delta_k^{(h)}. \quad (\text{V.3.14})$$

погрешности округления величин $\beta_{jkh}^{(k)}$ и c_k^h обозначим соответственно через $\delta_{jkh}^{(k)}$ и $\delta_k^{(h)}$. Пусть еще погрешности округления правых частей формул (VI.3.10-11), поделенных на $k!$, пусть соответственно $\tilde{\epsilon}_k^{(h)}$ и $\epsilon_k^{(h)}$. Положим, наконец, $\delta_{kh} = \max_{j \neq k} |\delta_{jkh}^{(k)}|$ и $\epsilon_{kh} = \max_{j \neq k} [|\tilde{\epsilon}_k^{(h)}| + |\beta_{jkh}^{(k)}| / k]$. Приняв во внимание, что (см. формулу (VI.2.11)

$$\sum_{j=0}^{2n-1} w_{je} \leq 2n = 1/h,$$

легко получить неравенство

$$\delta_{kh} \leq M(1 + \frac{1}{k})\delta_{k-1,h} + \epsilon_{kh}, \quad M = \max |f(t, \tau)|. \quad (\text{V.3.15})$$

Пусть вычисления производятся с такой точностью, что $\epsilon_{kh} \leq \epsilon = \text{const}$. Тогда можно получить (см. книгу автора [22]) следующую оценку (C_0 - некоторая постоянная, \bar{C} - постоянная Эйлера)

$$\delta_{kh} \leq \bar{C} \delta_{0h} k (C_0 M)^k - \frac{(\frac{3C_0 M}{2})^k - 1}{3C_0 M/2 - 1}. \quad (\text{VI.3.16})$$

Если ядро мало, так что $C_0 M < 2/3$, то δ_{kh} ограничена; при этом и λ невелико, то ограничена и погрешность второго ядра Фредгольма. Нетрудно убедиться, что при этих же условиях ограничена и погрешность определителя Фредгольма.

6°. Более оптимистическими оказываются вероятностные оценки погрешности округления. Допустим, что в суммах (VI.3.10-11), поделенных на k , ошибки округления частично сокращаются. При некоторых естественных предположениях (см. книгу автора [22]) это приводит к замене неравенства (VI.3.15) более слабым неравенством

$$\delta_{kh} \leq M h^d (1 + \frac{1}{k}) \delta_{k-1,h} + \epsilon_{kh}, \quad d = \text{const} > 0,$$

что в свою очередь приводит к оценке

$$\delta_{kh} \leq \bar{C} \delta_{0h} k (C_0 M h^d)^k + \frac{(\frac{3C_0 M h^d}{2})^k - 1}{3C_0 M h^d/2 - 1}, \quad (\text{VI.3.17})$$

которая выполняется с вероятностью, не меньшей, чем $1 - Ch^{3-2d}$,

$C = const$; таким образом, постоянную d можно выбрать по произволу из промежутка $(0, 3/2)$.

4. Сг денке к алгебраической системе

1°. Рассмотрим уравнение (VI.1.1) в следующих предположениях: а) $\mathcal{K} \in L_2([0, 1] \times [0, 1])$; б) единица не есть характеристическое число ядра $\mathcal{K}(t, \tau)$; в) $f \in L_2(0, 1)$. Пусть построено такое вырожденное ядро

$$\mathcal{K}^{(n)}(t, \tau) = \sum_{k=1}^n d_k(t) \beta_k(\tau), \quad (VI.4.1)$$

что

$$\int_0^1 \int_0^1 [\mathcal{K}(t, \tau) - \mathcal{K}^{(n)}(t, \tau)]^2 dt d\tau = r_n \xrightarrow{n \rightarrow \infty} 0. \quad (VI.4.2)$$

В качестве $\mathcal{K}^{(n)}(t, \tau)$ можно взять, например, конечноэлементную аппроксимацию ядра $\mathcal{K}(t, \tau)$, если только это последнее обладает необходимой гладкостью; отметим, что численные эксперименты подтверждали целесообразность такого выбора. Другие способы выбора ядра $\mathcal{K}^{(n)}(t, \tau)$ указаны в книге автора и Смолянского [1]. Будем считать, что свободный член уравнения $f(t)$ также заменен другим свободным членом $f^{(n)}(t)$, близким к $f(t)$ в норме $L_2(0, 1)$, так что $\|f - f^{(n)}\|_2 \xrightarrow{n \rightarrow \infty} 0$.

Таким образом, уравнение (VI.1.1) мы заменили новым уравнением

$$x^{(n)}(t) - \int_0^1 \mathcal{K}^{(n)}(t, \tau) x^{(n)}(\tau) d\tau = f^{(n)}(t); \quad (VI.4.3)$$

оценим погрешность аппроксимации такой замены. Пусть \mathcal{K} и $\mathcal{K}^{(n)}$ — операторы Фредгольма 2-го рода, соответственно равные $\mathcal{K}(t, \tau)$ и $\mathcal{K}^{(n)}(t, \tau)$. Положим $R = (I - \mathcal{K})^{-1}$, $R^{(n)} = (I - \mathcal{K}^{(n)})^{-1}$. Тогда $x_* = Rf$, $x_*^{(n)} = R^{(n)}f^{(n)}$ и

$$\|x_* - x_*^{(n)}\| = \|R(f - f^{(n)}) + (R - R^{(n)})f^{(n)}\| \leq$$

$$\leq \|R\| r_n' + \|R - R^{(n)}\| (\|f\| + r_n'); \quad r_n' = \|f - f^{(n)}\|. \quad (VI.4.4)$$

Далее, положим $\Delta^{(n)} = \mathcal{K} - \mathcal{K}^{(n)}$. Тогда $K^{(n)} = (I - \mathcal{K} + \Delta^{(n)})^{-1} = R(I + \Delta^{(n)}R)^{-1}$; отсюда $R - R^{(n)} = -R\Delta^{(n)}R(I + \Delta^{(n)}R)^{-1}$; если γ_n достаточно велико, то $\|\Delta^{(n)}R\| \leq \gamma_n \|R\| < 1$ и

$$|R - R^{(n)}| \leq \frac{r_n \|R\|^2}{1 - r_n \|R\|}$$

Подставив это в (VI.4.4), получим следующую оценку погрешности аппроксимации

$$\|x_n - x_n^{(n)}\| \leq \|R\| r_n + \frac{r_n \|R\|^2}{1 - r_n \|R\|} (\|f\| + r_n') \quad (\text{VI.4.5})$$

2°. Обычным способом решение уравнения (VI.4.3) сводится к решению некоторой линейной алгебраической системы ч порядка n ; при этом возникает погрешность искажения, которую можно оценить по формуле (I.4.14). Погрешности алгоритма и округления можно оценить по способу § 4 гл. I; оценки зависят от метода, которым система решается.

3°. Уравнение Фредгольма можно свести к линейной алгебраической системе по методу Бубнова - Галеркина. В этом случае погрешность аппроксимации можно оценить по формуле (II.8.1). Далее, если координатная система сильно минимальна в $L_2(0, 1)$, то процесс вычисления коэффициентов устойчив (в смысле второго определения) в последовательности (R_n, R_n) , а процесс вычисления приближенного решения устойчив в последовательности $(H^{(n)}, R_n)$; смысл обозначений очевиден.

4°. Еще один способ сведения к линейной алгебраической системе дает метод наименьших квадратов. Мы не будем здесь останавливаться на сущности метода - об этом достаточно полно сказано в книге автора [3], и сделаем только два замечания общего характера.

1) Пусть метод наименьших квадратов применен к уравнению $Ax = f$, где A - линейный оператор, действующий в некотором гильбертовом пространстве H , и что этот оператор ограниченно обратим. Если приближенное решение $x_n^{(n)}$ уже построено, то его погрешность оценивается очень просто [3]:

$$\|x_n - x_n^{(n)}\| \leq \|A\|^{-1} \|f - Ax_n^{(n)}\|. \quad (\text{VI.4.6})$$

Последняя формула дает суммарную погрешность приближенного решения.

2) Пусть еще оператор A ограничен в H . Тогда можно оценить погрешность аппроксимации метода наименьших квадратов по форму-

ле (IV.1.4).

Для уравнений Фредгольма обе упомянутые здесь формулы очевидно справедливы.

Сравнительно просто решается вопрос об устойчивости относительно искажения. В книге автора [3] приведены следующие простые факты. Если оператор A замкнут в гильбертовом пространстве H и имеет ограниченный обратный, если $f = \mathcal{D}(A^*)$, то уравнения $Ax = f$ и $A^*Ax = A^*f$ равносильны. Оператор A^*A положительно определен, его энергетическое пространство совпадает с множеством значений оператора A , а энергетическая норма $\|x\| = \|Ax\|$. Если A - ограниченный оператор, то, очевидно, эти последние нормы эквивалентны норме пространства H . Из сказанного сразу вытекает следующее утверждение, верное, если уравнение Фредгольма разрешимо единственным образом.

Для того чтобы процесс наименьших квадратов был устойчив относительно погрешностей искажения, необходимо и достаточно, чтобы координатная система была сильно минимална в $L_2(0, 1)$.

Система уравнений метода наименьших квадратов - алгебраическая линейная; для погрешностей аппроксимации и округления этой системы справедливо все сказанное по этому поводу в § 4 п. I.

Можно еще отметить, что если координатная система почти ортонормирована в $L_2(0, 1)$, то число обусловленности матрицы метода наименьших квадратов ограничено независимо от N .

Сказанное в пп. 3^о и 4^о настоящего параграфа полностью относится к уравнениям второго рода с вполне непрерывным оператором в сепарабельном гильбертовом пространстве.

5^о. К линейной алгебраической системе можно приближенно свести уравнение Фредгольма методом коллокации. В книге Пудгалова [3] доказывается, что оптимальная по порядку оценка погрешности аппроксимации достигается при использовании подходящим образом подобранных функций в качестве координатных функций. Если ядро и свободный член уравнения (VI.1.1) принадлежат к классу C^k , а их модули непрерывны относительно одной и той же функции $\omega(\eta)$, то погрешность аппроксимации коллокационным методом имеет оценку $O(n^{-k} \omega(n^{-1}))$. В той же книге приведены и некоторые другие результаты, близкие к только что упомянутому. По поводу точности процесса коллокации относительно погрешностей искажения можно повторить то, что было

сказано в общем случае в § 9 гл. III. В частности, если выбрана последовательность координатных функций $y_k(t)$, $1 \leq k < \infty$, и $|y_k'(t)| \leq Ck^\alpha$, $\alpha < 1$, то процесс коллокации неустойчив.

При фиксированном порядке приближения и при заданных координатных функциях и узлах коллокации мы имеем дело с вполне определенной алгебраической системой. Для этой системы можно оценить погрешности искажения, алгоритма и округления так, как это описано в § 4 гл. I.

§ 5. Сингулярные интегральные уравнения. Метод наименьших квадратов

1°. Будем рассматривать одномерное сингулярное интегральное уравнение вида

$$Ax = -a(t)x(t) + t(t)(Sx)(t) + (Tx)(t) - f(t), \quad (VI.5.1)$$

где S - оператор Коши

$$(Sx)(t) = \frac{1}{\pi i} \int_{\Gamma} \frac{x(\tau) d\tau}{\tau - t}, \quad (VI.5.2)$$

в следующих предположениях: Γ - замкнутый плоский контур класса $C^{(1, \alpha)}$, ограничивающий некоторую область, односвязную или многосвязную, $a(t)$ и $t(t)$ - функции (скалярные или матричные), непрерывные на Γ , T - оператор, вполне непрерывный в $L_2(\Gamma)$, индекс уравнения (VI.5.1) равен нулю и это уравнение имеет не более одного решения, так что оно разрешимо при любой (соответственно скалярной или векторной) функции $f(t)$. При перечисленных условиях операторы A и A^{-1} ограничен в $L_2(\Gamma)$.

2°. Уравнение (VI.5.1) можно решать по методу наименьших квадратов, и к этому уравнению в полной мере относится то, что было сказано в п. 4° предшествующего параграфа.

Отметим некоторые координатные системы.

1) Пусть $\Gamma = \bigcup_{j=0}^q \Gamma_j$, где Γ_j - замкнутые кривые; предположим, что Γ есть граница некоторой конечной области D и что Γ_0 ограничивает область D извне. Зная Γ_j , $1 \leq j \leq q$, выберем точку t_j . Обозначим через $\rho_0, \rho_1, \dots, \rho_q$ числа, удовлетворяющие неравенствам

$$0 < \rho_0 < \min_{t \in \Gamma} r, \quad \rho_j > \max_{t \in \Gamma_j} |t - t_j|, \quad 1 \leq j \leq q$$

Тогда система функций

$$\left(\frac{t}{\beta_0}\right)^{n-1}, \left(\frac{t}{\beta_0}\right)^{n-1} \cdot \left(\frac{\beta_1}{t-\beta_1}\right)^n, \left(\frac{\beta_2}{t-\beta_2}\right)^n, \quad n=1, 2, \dots, \quad 1 \leq j \leq q$$

слабо минимальна в $L_2(\Gamma)$ (см. книгу автора [3]).

2) Пусть Γ - вещественная ось, и пусть вполне непрерывный оператор T - сглаживающий, так что $T: W_2^{(k)}(R_1) \rightarrow W_2^{(k+1)}(R_1)$ holds, и в качестве координатных взять функции

$$\varphi_k^{(1)}(t) = \sqrt{\frac{t}{1+t^2}} T_k\left(\frac{2t}{1+t^2}\right), \quad k=0, 1, 2, \dots,$$

$$\varphi_k^{(2)}(t) = \frac{\sqrt{2}(1-t^2)}{(1+t^2)^{3/2}} U_{k-1}\left(\frac{2t}{1+t^2}\right), \quad k=1, 2, \dots,$$

где T_k и U_k - полиномы Чебышева соответственно первого и второго рода, то общую оценку (L.I.5) можно уточнить: при $f \in W_2^{(k)}$, $a(t), b(t) \in W_2^{(k)}(R_1)$ будет

$$\|x_n - x_n^{(n)}\|_{2, \Omega_n} = O(n^{-3}).$$

При этом процесс наименьших квадратов устойчив относительно граничных искажения, а число обусловленности матрицы метода наименьших квадратов ограничено (см. книгу автора и Преслора [1]).

3°. Для простоты допустим, что Γ - окружность $|z|=1$ и что уравнение (VI.5.1) - скалярное. Допустим, что решение уравнения (VI.5.1) $x_n \in W_2^{(2\delta)}$; для этого достаточно, например, чтобы вполне непрерывный оператор T был сглаживающим (или чтобы он действовал из $L_2(\Gamma)$ в $W_2^{(2\delta)}(\Gamma)$), а $f(t), a(t), b(t) \in W_2^{(2\delta)}$. В качестве координатных возьмем функции, которые получаются кусочно-полиномиальными, степени $2\delta-1$, функций МКЭ (см. книгу автора [22], гл. II, § 6). В этом случае, как нетрудно доказать, погрешность аппроксимации имеет оценку $O(n^{-2\delta})$. Повторяя, некоторыми естественными упрощениями, рассуждения книги автора [22], гл. IX, § 2, легко доказать, что при выборе координатных функций, указанных в данном пункте, справедливо следующее утверждение: процесс наименьших квадратов устойчив относительно граничных искажения в последовательности гур-пространств (X_n, \bar{Y}_n) , где X_n и \bar{Y}_n - n -мерные пространства с нормами

$$\forall a \in R_n, \|a\|_{\tilde{X}_n} = h^{1/2} \|a\|_{R_n}, \|\tau\|_{\tilde{X}_n} = h^{1/2} \|a\|_{R_n}.$$

Точно так же доказывается, что процесс построения приближенного решения по методу наименьших квадратов устойчив в последовательности (H_n, Y_n) , где H_n есть n -мерное подпространство пространства $L_2(\Gamma)$, натянутое на соответствующие координатные функции. Наконец, рассуждения § 5 гл. IX книги автора [22] показывают, что число обусловленности матрицы метода наименьших квадратов при том же выборе координатных функций с n связано независимо от n .

Распространение результатов данного пункта на системы одномерных сингулярных уравнений не встречает затруднений.

4°. Отметим результаты статьи [12] Джикхариани, в которой исследуются одномерные сингулярные интегральные уравнения на разомкнутом контуре. В цитируемой статье Джикхариани ограничивается случаем уравнений

$$ax(t) - \frac{b}{\pi i} \int_{-1}^1 \frac{x(\tau) d\tau}{\tau-t} + \frac{1}{\pi} \int_{-1}^1 K(t,\tau)x(\tau) d\tau = f(t), \quad (VI.5.3)$$

где a, b - постоянные, $a^2 + b^2 > 0$. Уравнение рассматривается в пространстве $H = L_2(-1, +1; \rho)$; вес $\rho(t)$ выбирается в зависимости от индекса оператора (VI.5.3). Интегральный оператор

$$\int_{-1}^1 K(t,\tau)x(\tau) d\tau$$

предполагается вполне непрерывным в H ; предполагается также, что уравнение (VI.5.3) имеет в H единственное решение. Приближенное решение этого уравнения отыскается по методу Бубнова - Галеркина, в качестве координатных функций выбираются полиномы Чебышева или Эрмита. Доказывается, что алгебраическая система Бубнова - Галеркина достаточно высокого порядка однозначно разрешима и что приближенные решения сходятся к точному решению в H . Утверждается устойчивость процесса.

5°. Перейдем к рассмотрению многомерных сингулярных уравнений. Рассмотрим сингулярное (скалярное или векторное, безразлично) интегральное уравнение

$$Ax := a(t)\sigma(t) + \int_{\Gamma} K(t,\tau)x(\tau) d\tau + (\tau x)(t) = f(t). \quad (VI.5.4)$$

Здесь Γ - достаточно гладкое многообразие m измерений без края, $\mathcal{K}(t, \tau)$ - сингулярное ядро, T - оператор, вполне непрерывный в том банаховом пространстве, в котором сингулярный оператор ограничен; для определенности будем говорить о пространстве $L_2(\Gamma)$. Примем следующие допущения: 1) символ $\Phi(t, \tau)$ уравнения (VI.5.1) - достаточно гладкая функция своих аргументов; 2) левая часть уравнения равна нулю; 3) уравнение (VI.5.4) имеет в $L_2(\Gamma)$ не более одного решения - тогда от разрешимости любой правой части из $L_2(\Gamma)$ операторы A и A' ограничены. Повторяя многократно использованные выше рассуждения, мы приходим к следующим результатам.

Уравнение (VI.5.4) допускает применение метода наименьших квадратов; погрешность аппроксимации оценивается по формуле (IV.1.4). Если координатная система сильно минимальна в L_2 , то соответствующий вычислительный процесс устойчив относительно погрешностей измерения; если эта система еще и почти ортонормирована, то число обусловленности алгебраической системы методов наименьших квадратов ограничено независимо от порядка этой системы. Таким образом, если эту систему решить по методу простых итераций, то эти итерации сходятся, как прогрессия с фиксированным множителем, меньшей единицы, и погрешности округления возрастания порядка системы остаются ограниченными.

В статьях автора и Раевой [1] и автора [25] приведены примеры координатных систем, для которых общая оценка (IV.1.4) может быть уточнена.

§ 6. Методы механических квадратур и колlocation для одномерных сингулярных уравнений.

1°. Приложениям методов механических квадратур к одномерным сингулярным интегральным уравнениям посвящен ряд работ Габдуллаева и некоторых других авторов. Подробно полная сложность приведена в статье Габдуллаева и Душкова [1]; мы ограничимся здесь кратким изложением статьи Габдуллаева [2].

Рассмотрим уравнение (VI.5.1), в котором вполне непрерывный оператор T - интегральный:

$$(Tx)(t) = \int_{\Gamma} \mathcal{K}(t, \tau) x(\tau) d\tau. \quad (VI.5.1)$$

Примем, что Γ - окружность $|z| = 1$, и что функции $a(t)$, $b(t)$, $f(t)$, $\mathcal{K}(t, \tau)$ принадлежат к классу $C^{(r, d)}$, где

$\nu > 0$, $0 < \lambda < 1$. Допустим далее, что индекс рассматриваемого уравнения равен нулю и что это уравнение имеет не более одного решения. Приближенное решение будем искать в виде тригонометрического многочлена

$$x^{(n)}(t) = \sum_{k=-n}^n a_k t^k. \quad (VI.6.2)$$

Из элементарных свойств оператора Коши вытекает, что

$$(Ax^{(n)})(t) =$$

$$[a(t)+b(t)] \sum_{k=0}^n a_k t^k + [a(t)-b(t)] \sum_{k=-1}^{-n} a_k t^k + \int_{\Gamma} \gamma(t, \tau) \sum_{k=-n}^n a_k \tau^k d\tau.$$

Изменим здесь интеграл его приближенным значением по формуле прямоугольников, соответствующей разбиению окружности Γ на $2n+1$ равных дуг с точками деления $t_j = \exp\{2\pi i j / (2n+1)\}$, и потребуем, чтобы полученная таким образом функция от t совпала с $x^{(n)}(t)$ в точках t_j . Это приведет нас к системе линейных алгебраических уравнений

$$\begin{aligned} & [a(t_j)+b(t_j)] \sum_{k=0}^n a_k t_j^k + [a(t)-b(t)] \sum_{k=-1}^{-n} a_k t_j^k + \\ & + \sum_{k=-n}^n a_k \sum_{\rho=0}^{2n} \mathcal{K}(t_j, t_{\rho}) [t_{\rho}^k - t_{\rho}^{-k}] = f(t_j); \quad j=0, 1, \dots, 2n. \end{aligned} \quad (VI.6.3)$$

В статье Габдулхаева [...] эта система трактуется как с.з.м. метода механических квадратур для уравнения (VI.5.1). Может быть, стоит отметить, что метод, примененный Габдулхаевым представляет собой некоторый синтез методов коллокации и механических квадратур: он переходит в метод коллокации при $\gamma(t, \tau) \equiv 0$ и в метод механических квадратур при $b(t) \equiv 0$.

Основной результат цитируемой статьи содержится в следующей теореме:

Теорема VI.6.1. Пусть β - любое число, ε - любое в пределах $0 < \beta < \varepsilon$. Существуют такие положительные постоянные A_1, B_1, A_2, B_2 , что если

$$(A_1 \ln n + B_1) n^{-\varepsilon} < 1, \quad (VI.6.4)$$

то система (У1.6.3) имеет единственное решение. Истинность проксимации приближенного решения (У1.6.2) оценивается по формуле

$$\|x_* - x_*^{(n)}\|_{H_\beta} \leq (A_2 \ln n + B_2) n^{-\tau+\beta}; \quad (У1.6.4)$$

H_β - пространство функций, заданных на Γ , с нормой

$$\|x\|_{H_\beta} = \max_{t \in \Gamma} |x(t)| + \sup_{t', t'' \in \Gamma} \frac{|x(t') - x(t'')|}{|t' - t''|^\beta}.$$

Доказательство этой теоремы использует общую теорию приближенных методов, построенную Канторовичем (см. Канторович и Акилов [1]).

2°. Рассмотрим теперь уравнение (У1.5.1) в следующих предположениях: $T=0$; $a(t)$, $b(t)$ непрерывны на Γ ; индекс влечения равен нулю. Заметим, что результаты статьи Габдулхаева [2] здесь неприменимы, потому что $\tau = d = 0$. Данный случай рассмотрен в статье Пресдорфа и Шмидта [1]. В качестве узлов, локализации брать точки $t_k^{(n)} = \exp(2\pi i k/n)$, а в качестве координатных функций - кусочно линейные функции, широко используемые в МКЭ и определяемые формулой

$$y_k^{(n)}(t) = \begin{cases} \frac{t - t_{k-1}}{t_k - t_{k-1}}; & \frac{2\pi(k-1)}{n} \leq \arg t < \frac{2\pi k}{n}; \\ \frac{t_{k+1} - t}{t_{k+1} - t_k}; & \frac{2\pi k}{n} \leq \arg t < \frac{2\pi(k+1)}{n}; \\ 0 & \text{в остальных случаях;} \end{cases} \quad (У1.6.5)$$

Приближенное решение имеет вид

$$x^{(n)}(t) = \sum_{k=0}^{n-1} a_k^{(n)} y_k^{(n)}(t); \quad (У1.6.6)$$

коэффициенты $a_k^{(n)}$ определяются из систем

$$a_j^{(n)} a_j^{(n)} + b(t_j) \sum_{k=0}^{n-1} a_k^{(n)} (S y_k^{(n)})(t_j) = f(t_j); \quad j=1, 2, \dots, n. \quad (У1.6.7)$$

Основные результаты статьи Пресдорфа и Шмидта дают следующие теоремы:

Теорема У1.6.2. Пусть коэффициенты $a(t)$, $b(t)$ непрерывны, а решение $x_* \in L_2(\Gamma)$. Пусть выполнено условие

$$\forall \lambda \in [-1, 1], \forall t \in \Gamma; a(t) + \lambda b(t) \neq 0. \quad (VI.6.9)$$

Тогда при достаточно больших n система (VI.6.8) разрешима единственнм образом и $\|x_* - x_*^{(n)}\|_{2,\Gamma} \xrightarrow{n \rightarrow \infty} 0$.

Теорема VI.6.3. Если $a(t), b(t) \in H_\mu, 0 < \mu < 1$, и $f \in H_\nu, 0 < \nu < 1$, то справедлива оценка

$$\|x_* - x_*^{(n)}\|_{2,\Gamma} \leq C n^{-\sigma} \ln n, \quad \sigma = \min(\mu, \nu).$$

Теорема VI.6.4. Пусть $a(t), b(t) \in H_\mu, 1/2 < \mu < 1$.

Если при любой непрерывной функции $f(t)$ система (VI.6.8) разрешима при достаточно большом n , и $\|x_* - x_*^{(n)}\|_{2,\Gamma} \xrightarrow{n \rightarrow \infty} 0$,

то справедливо неравенство (VI.6.9).

3°. При фиксированном n можно получить оценки для погрешностей искоженной, алгоритма и округления так, как об этом сказано в § 4 гл. I.

4°. Дальнейшие [14] изменили метод коллокации к уравнениям вида (VI.5.2) при ограничениях, перечисленных в п. 4° § 5 данной главы.

Доказывается, что если за узлы коллокации взять корни полиномов Лежандра, выбранных подходящим образом, то алгебраически системы метода коллокации, порядок которых достаточно велик, однозначно разрешимы, и приближенные решения сходятся к точным решениям соответствующей матрице.

ГЛАВА II

Нелинейные вычислительные процессы

§ 1. О г-решности аппроксимации метода Рунге

В п. 1.0 - 4.0 настоящего параграфа мы коротко изложим постановку задачи о минимуме функционала и метод Рунге для ее решения; подробно эти вещи изложены в книге автора [8], гл. IX. В п. 5 будет рассмотрена полнота аппроксимации процесса Рунге

1.0. Пусть B - рефлексивное банахово пространство и F - функционал, определенная на линейном множестве $D(F)$, плотном в B . Ограничимся случаем, когда функционал F непрерывно дифференцируем на любой конечномерной гиперплоскости, лежащей в B . Конечномерной гиперплоскостью в банаховом пространстве B мы назовем совокупность элементов вида

$$x_0 + \sum_{k=1}^n a_k x_k, \quad (УП.1.1)$$

где n - фиксированное натуральное число, называемое размерностью гиперплоскости; $a_k, 1 \leq k \leq n$ - произвольные числа, $x_k, 0 \leq k \leq n$ - фиксированные элементы пространства B . Говоря, что F обладает той или иной гладкостью на гиперплоскости (УП.1.1), мы имеем под этим, что такой гладкостью обладает функция переменных (a_1, a_2, \dots, a_n) , равная $F(x_0 + a_1 x_1 + \dots + a_n x_n)$.

Если функционал F имеет локальный минимум в точке $x_0 \in D(F)$, то

$$(\text{grad } F)(x_0) = 0. \quad (УП.1.2)$$

Мы будем обозначать $\text{grad } F = P$; очевидно, P действует из B в B^* .

Функционал F называется выпуклым, если

$$\forall x, y \in D(F), \forall \lambda \in [0, 1], F(\lambda x + (1-\lambda)y) \leq \lambda F(x) + (1-\lambda)F(y); \quad (УП.1.3)$$

он называется существенно выпуклым, если в соотношении (УП.1.3) равенства имеет место только при $x = y$. Для функции F называется точкой стационарной, если $F(x) \rightarrow +\infty$ тогда и только тогда, когда $\|x\| \rightarrow \infty$. Функционал F называется полунепрерывным снизу (с.р.) в точке x_0 , если $\liminf_{x \rightarrow x_0} F(x) \geq F(x_0)$ (соответств. же, $\limsup_{x \rightarrow x_0} F(x) \leq F(x_0)$). Тот же функционал называется слабо полунепрерывным снизу (с.р.), если упомянутые

соотношения выполняются каждый раз, когда $x \rightarrow x_0$; символ \rightarrow означает слабую сходимость.

Теорема VII.1.1. Если функционал \mathcal{F} возрастающий, существенно вынужденный и слабо полунепрерывный снизу, то он ограничен снизу и его нижняя грань достигается в единственной точке, к которой слабо сходится любая минимизирующая последовательность.

Эта теорема доказана в работе Гельмана [1].

2°. В данном пункте излагаются результаты работы Бирман [1]

Пусть линейный оператор A действует из рефлексивного пространства B в сопряженное пространство B^* , тогда сопряженный оператор A^* также действует из B в B^* . Будем говорить, что оператор $A : B \rightarrow B^*$ симметричен, если $\mathcal{D}(A) \subset \mathcal{D}(A^*)$ и $Ax = A^*x$, $x \in \mathcal{D}(A)$; 2) самосопряжен, если он симметричен и $\mathcal{D}(A) = \mathcal{D}(A^*)$; 3) положителен, если он симметричен и $(Ax, x) > 0$ при $x \neq 0$; 4) положительно определен, если он симметричен и $(Ax, x) \geq \gamma^2 \|x\|^2$, где γ - положительная постоянная. Если A - положительный (в частности, положительно определенный) оператор, то с ним можно связать энергетическое пространство B_A так же, как в случае, когда пространство B само гильбертово (см. книгу автора [3]). Именно, за B_A мы принимаем гильбертово пространство, получаемое замыканием множества $\mathcal{D}(A)$ в норме, порожденной скалярным произведением $[x, y]_A = [x, y] = (Ax, y)$; здесь (Ax, y) есть значение функционала $Ax \in B^*$ на элементе $y \in B$. Норму в B_A будем обозначать через $|\cdot|$, так что $|x| = \sqrt{[x, x]}$.

Если A положительно определенный оператор, то B_A конечномерно вкладывается в B ; для положительного, но не положительно определенного оператора это утверждение неверно.

3°. Пусть оператор $P = \text{grad } \mathcal{F}$ имеет производную Гато P'_u со следующими свойствами: а) эта производная существует при любых $u \in \mathcal{D}(P)$; б) $\mathcal{D}(P'_u) \supset \mathcal{D}(P)$; в) оператор P'_u - положительно определенный при любых $u \in \mathcal{D}(P)$, так что

$$\forall x \in \mathcal{D}(P'_u), (P'_u x, x) \geq \gamma^2 \|x\|^2, \quad (\text{VII.1.4})$$

где величина γ не зависит от u . Доказывается, что в этих условиях функционал \mathcal{F} ограничен снизу, и любая минимизирующая последовательность сходится в метрике B к пределу, который не зависит от выбора минимизирующей последовательности. Это означает

будем рассматривать как обобщенное решение задачи с минимумом функционала \mathcal{F} .

Представляет интерес случай, когда P_u удовлетворяет не только неравенству (УП.1.1), но и более сильному неравенству

$$\forall x \in \mathcal{D}(P'_u), (P'_u x, u) \geq d^2 (P'_u x, x), d = \text{const} > 0, \text{ (УП.1.2)}$$

где u_0 - некоторый фиксированный элемент из $\mathcal{D}(P'_u)$. Обозначим через U_0 энергетическое пространство оператора P'_u и через $(\cdot, \cdot)_0$ скалярное произведение и норму в этом пространстве. Можно доказать, что в этом случае обобщенное решение принадлежит пространству B_0 .

4^o. Коротко опишем метод Рунге для задачи о минимуме функционала \mathcal{F} . Напомним, что частным случаем метода Рунге является МКЭ. Зададим последовательность конечных подпространств координатных элементов

$$\{y_{n1}, y_{n2}, \dots, y_{nN}\}; n=1, 2, \dots; N=N(n), \text{ (УП.1.6)}$$

подчиненных условиям а) $y_{nk} \in \mathcal{D}(\mathcal{F}) \cap B_0$; б) при данном n элементы y_{nk} линейно независимы; в) последовательности (УП.1.6) полна в B_0 . Обозначим через B_n подпространство пространства B_0 с базисом $\{y_{nk}\}$, $k=1, 2, \dots, N$, и будем искать минимум \mathcal{F} на B_n . Решение $x_*^{(n)}$ этой задачи - приближенное решение данной задачи по Рунге - существует и единственно, если \mathcal{F} удовлетворяет условиям пп. 2^o, 3^o; оно приближенное решение имеет вид

$$x_*^{(n)} = \sum_{k=1}^N a_k^{(n)} y_{nk}, \text{ (УП.1.7)}$$

где $a_k^{(n)}$ - численные коэффициенты, определяемые из линейной системы Рунге

$$\frac{\partial}{\partial a_j^{(n)}} \mathcal{F} \left(\sum_{k=1}^N a_k^{(n)} y_{nk} \right) = 0, 1 \leq j \leq N. \text{ (УП.1.8)}$$

Будем рассматривать только то решение системы (УП.1.8), которое доставляет величине $\mathcal{F} \left(\sum_{k=1}^N a_k^{(n)} y_{nk} \right)$ абсолютный минимум (эта оговорка не нужна, если функционал \mathcal{F} существенно вогнутый - в этом случае система (УП.1.8) имеет единственное решение).

Если последовательность подпространств B_n расширяющаяся, то (см. книгу автора [8]) последовательность $\{x_*^{(n)}\}$ приближенных

решений по Ритцу - минимизирующая для функционала \mathcal{F} , если этот последний - возрастающий и полунепрерывный сверху. Нетрудно показать, что это утверждение верно. Тогда, когда подпространства B_n не расширяются (так будет, например, в случае МКЭ), но их последовательность полна в B_0 . в цитированной выше книге автора [8], § 70, доказывается (без предположения о том, что пространство B_n расширяется), что из любой последовательности приближенных решений по Ритцу можно выделить минимизирующую подпоследовательность. Допустим, что последовательность всех приближенных решений по Ритцу не минимизирующая. Тогда существует такое число $\varepsilon_0 > 0$ и такая подпоследовательность $\{x_n^{(n)}\}$ приближенных решений по Ритцу, что $\mathcal{F}(x_n^{(n)}) > d + \varepsilon_0$, $d = \inf \mathcal{F}(x)$.

Но из этой подпоследовательности можно выделить минимизирующую последовательность, что невозможно. Полученное противоречие доказывает наше утверждение.

Из доказанного утверждения следует, что о сходимости метода Ритца можно поговорить все сказанное в п.3^о о сходимости минимизирующей последовательности.

5^о. Для некоторого класса функционалов можно получить оценку погрешности аппроксимации процесса Ритца. Допустим, что проволочная P'_u оператора $P = \text{grad } \mathcal{F}$ удовлетворяет, кроме неравенств (УП. I.4) к (УП. I.5), еще и следующую

$$\forall x \in \mathcal{D}(P'_u), (P'_u x, x) \leq \beta^2 (P'_u x, \tau), \quad \beta = \alpha nst > 0. \quad (\text{УП. I.9})$$

Воспользуемся равенством

$$\mathcal{F}(x_n + y) - \mathcal{F}(x_n) = \int_0^1 \tau^{-1} d\tau \int_0^1 (P_{x_n + \tau y} \tau y, \tau y) dt.$$

Изменив равенство (УП. I.9) и положив $x_n + y = x$, получим

$$\mathcal{F}(x) - \mathcal{F}(x_n) = \beta^2 \int_0^1 \tau^{-1} d\tau \int_0^1 (P'_u(\tau(x-x_n)), \tau(x-x_n)) dt = \frac{1}{2} \beta^2 \|x - x_n\|_0^2. \quad (\text{УП. I.10})$$

Пусть $y^{(n)}$ - элемент пространства B_n , на котором достигается наилучшее приближение элемента x_* :

$$\varepsilon_n(x_*) = \inf_{x^{(n)} \in B_n} \|x_* - x^{(n)}\|_0 = \|x_* - y^{(n)}\|_0.$$

Подставив в (УП. I.10) $x = y^{(n)}$, находим

$$F(y^{(n)}) - F(x_*) \leq \frac{1}{2} \beta^2 |y^{(n)} - x_*|^2 = \frac{1}{2} \beta^2 \varepsilon_n^2(x_*). \quad (\text{УП. I. 1})$$

Если $x_*^{(n)}$ - приближенное решение по Рунге, то $F(x_*^{(n)}) - F(y^{(n)})$ по неравенствам (УП. I. 5) и (УП. I. 9) получаем

$$\frac{1}{2} \alpha^2 |x_* - x_*^{(n)}|^2 \leq F(x_*^{(n)}) - F(x_*) \leq F(y^{(n)}) - F(x_*) \leq \frac{1}{2} \beta^2 \varepsilon_n^2(x_*),$$

что дает следующую оценку погрешности аппроксимации:

$$|x_* - x_*^{(n)}| \leq \frac{\beta}{\alpha} \varepsilon_n(x_*). \quad (\text{УП. I. 1})$$

Если пространство B - функциональное, то оценка порядка величины $\varepsilon_n(x_*)$ как функции от n существенно зависит от гладкости решения x_* . К этому тов. у ом., например, книги Мо [1], Леценконой и Уральцевой [1]. Для метода конечных элементов оценки величины $\varepsilon_n(x_*)$ в зависимости от гладкости функции у ом. книги автора [17, 22], Обена [1], Странга и Экса [1], Сьале [1].

§ 2. Погрешность нахождения и устойчивость свободного вычислительного процесса

Г^с. Пусть даны две последовательности банаховых пространств X_n и Y_n , $n = 0, 1, 2, \dots$, и пусть при любом n оператор (и вообще говоря, нелинейный) A_n действует из X_n в Y_n , определен на всем пространстве X_n (это требование можно ослабить), удовлетворяет равенству $A_n(0) = 0$, имеет обратный оператор A_n^{-1} , который в свою очередь определен на всем пространстве Y_n и удовлетворяет следующему условию Липшица: если $\|y_1\|, \|y_2\| \leq t$,

$$\|A_n^{-1}(y_1) - A_n^{-1}(y_2)\| \leq C_n(t) \|y_1 - y_2\|, \quad (\text{УП. 2. 1})$$

где $C_n(t)$ не зависит от y_1 и y_2 .

Здесь и ниже, там, где это не может вызвать недоразумения мы не ставим обозначения пространства при члене нормы.

Рассмотрим свободный вычислительный процесс, состоящий в решении последовательности независимых уравнений

$$A_n(x^{(n)}) = f^{(n)}; \quad n = 0, 1, 2, \dots, \quad (\text{УП. 2. 2})$$

и соответствующий итерационный процесс

$$A_n(z^{(n)}) + \Gamma_n(z^{(n)}) = f^{(n)} + \delta^{(n)}. \quad (\text{УП.2.3})$$

Допустим, что искажения Γ_n и $\delta^{(n)}$ зависят от параметра $\delta > 0$ так, что

$$\|\delta^{(n)}\| \leq \delta \quad (\text{УП.2.4})$$

и существуют числовые функции $g_n(t, \delta)$ и $\psi_n(t, \delta)$ положительных числовых переменных t и δ со следующими свойствами:

$\lim_{\delta \rightarrow 0} g_n(t, \delta) = \lim_{\delta \rightarrow 0} \psi_n(t, \delta) = 0$; если $\alpha, \alpha', \alpha'' \in X_n$ и нормы этих элементов не превосходят t , то

$$\|\Gamma_n(\alpha)\| \leq g_n(t, \delta), \quad (\text{УП.2.5})$$

$$\|\Gamma_n(\alpha') - \Gamma_n(\alpha'')\| \leq \psi_n(t, \delta) \|\alpha' - \alpha''\|. \quad (\text{УП.2.6})$$

Уравнение (УП.2.3) равносильно следующему

$$z^{(n)} = A_n^{-1}(-\Gamma_n(z^{(n)}) + f^{(n)} + \delta^{(n)}). \quad (\text{УП.2.7})$$

Запишем Π и обозначим $\|f^{(n)}\| = T_n$. Докажем, что оператор в правой части уравнения (УП.2.7) есть оператор сжатия в шаре $S_n = \{z^{(n)} \in X_n : \|z^{(n)}\| \leq 2T_n C_n(T_n)\}$, если только δ достаточно мало. Имеем

$$\|-\Gamma_n(z^{(n)}) + f^{(n)} + \delta^{(n)}\| \leq g_n(2T_n C_n(T_n), \delta) + T_n + \delta,$$

что не превосходит $2T_n$ при достаточно малом δ . Теперь по неравенству (УП.2.1) $\|A_n^{-1}(-\Gamma_n(z^{(n)}) + f^{(n)} + \delta^{(n)})\| \leq 2T_n C_n(T_n)$

и оператор (УП.2.7) отображает шар S_n в себя. Далее, если $z^{(n)} \in S_n$, то $\|z^{(n)}\| \leq 2T_n C_n(T_n)$, и по неравенству (УП.2.1) и (УП.2.6)

$$\|A_n^{-1}(-\Gamma_n(z^{(n)}) + f^{(n)} + \delta^{(n)}) - A_n^{-1}(-\Gamma_n(z^{(n)}) + f^{(n)} + \delta^{(n)})\| \leq$$

$$\leq C_n(2T_n C_n(T_n)) \|\Gamma_n(z^{(n)}) - \Gamma_n(z^{(n)})\| \leq$$

$$\leq C_n(2T_n C_n(T_n)) \psi_n(2T_n C_n(T_n), \delta) \|z^{(n)} - z^{(n)}\|;$$

если δ достаточно мало, то множитель при норме $\|x^{(n)} - x^{(n)}\|$ меньше единицы, и наше утверждение доказано. В силу принципа сужения отображений уравнение (УП.2.3) имеет в шаре S_n ровно одно решение. Если $x_*^{(n)}$ — это решение, то

$$\|x_*^{(n)} - x_*^{(n)}\| = \|A_n^{-1}(\Gamma_n(x_*^{(n)})\{f + \delta^{(n)}\} - A_n(f^{(n)}))\| <$$

$$< C_n(2T_n) \|\Gamma_n(x_*^{(n)}) - \delta^{(n)}\| \leq C_n(2T_n) [y_n(2T_n C_n(T_n), \delta) + \delta]. \quad (\text{УП.2.4})$$

Формула (УП.2.3) дает оценку погрешности искажения для процесса (УП.2.2) в довольно общих условиях; при фиксированном n эта оценка стремится к нулю вместе с δ . Представляет интерес случай, когда указанное стремление происходит равномерно относительно n . С этим мы займемся в следующем пункте.

2°. В книге автора [8] дано некоторое определение устойчивости стохастического нелинейного вычислительного процесса. Получена теорема о достаточных условиях устойчивости процесса Ритца для некоторого класса нелинейных уравнений. Несколько более общее определение устойчивости нелинейного стохастического процесса дал Таккер [1]. Смолден в своей диссертации [1] доказал некоторые теоремы об устойчивости по Таккеру и применял их для исследования устойчивости МКЭ.

Мы дадим здесь другое определение устойчивости естественным образом обобщающее определение гл. III.

Допустим, что в некотором шаре $\|x^{(n)}\| \leq \bar{t}$ уравнение (УП.2.1) однозначно разрешимо при любом n . Относительно искаженного процесса (УП.2.2) допустим, что справедливы неравенства (УП.2.4) — (5) причем функции φ и ψ_n не зависят от n ; таким образом если $x, x', x'' - X$ и нормы этих элементов не превосходят числа \bar{t} , то

$$|\Gamma_n(x)| \leq \varphi(t, \delta), \quad (\text{УП.2.9})$$

$$\| \Gamma_n(x') - \Gamma_n(x) \| \leq \psi(t, \delta) \|x' - x\|, \quad (\text{УП.2.10})$$

$$\lim_{\delta \rightarrow 0} \varphi(t, \delta) = \lim_{\delta \rightarrow 0} \psi(t, \delta) = 0. \quad (\text{УП.2.11})$$

Назовем процесс (УП.2) устойчивым в последовательности пар процессов (X_n, Y_n) , если для любого $\varepsilon > 0$ существует такое

$\delta > 0$, ч при выполнении неравенств (УП.2.4), (УП.2.9-II) уравнение (УП.2.3) однозначно разрешимо в некоторой окрестности точки $X_*^{(n)}$, причем радиус этой окрестности не зависит от n , и справедливо неравенство

$$\|x^{(n)} - X_*^{(n)}\| < \varepsilon. \quad (\text{УП.2.12})$$

Теорема УП.2.1 Пусть $\|f^{(n)}\| \leq T = \text{const}$, $A_n(0) = 0$, оператор A_n определен в шаре $\|f^{(n)}\| \leq T$ и удовлетворяет в этом шаре условию Липшица

$$\|A_n^{-1}(f_1) - A_n^{-1}(f_2)\| \leq C \|f_1 - f_2\| \quad (\text{УП.2.13})$$

с постоянной C , которая не зависит от n . Пусть еще искажения $\delta^{(n)}$ и Γ_n удовлетворяют неравенствам (УП.2.4), (УП.2.9-II). Тогда процесс (УП.2.2) устойчив в последовательности (A_n, Y_n) .

Для доказательства достаточно заметить, что в условиях теоремы УП.2.1 радиус шара S_n , в котором разрешимо искаженное уравнение (УП.2.3), не зависит от n , и что правая часть формулы (УП.2.8) также не зависит от n и стремится к нулю вместе с δ .

Замечание УП.2.1. Условие $\|f^{(n)}\| \leq T$ можно сбросить, если оператор A_n^{-1} определен на всем пространстве Y_n и удовлетворяет условию Липшица (УП.2.13), в котором C не зависит от n .

§ 3. Об устойчивости процессов Рунге и конечных элементов

1°. В книге автора [3] § 76, доказано, что классический процесс Рунге для задачи $F(x) = \min$ устойчив, если: 1) оператор $P = \text{grad } F$ имеет производную Гато P'_u , удовлетворяющую неравенствам (УП.1.4-5); 2) координатная система сильно минимальна в энергетическом пространстве оператора P'_u . Точность доказательства проведена в предположении, что искажения Γ_n удовлетворяют требованиям, несколько более ограничительным, чем требования (УП.2.9-10). Однако, легко провести доказательство, предполагая выполненными только последние требования. Мы не будем останавливаться на этом подробнее; ниже мы докажем более общую теорему, аналогичную теореме Ш.З. для линейных задач.

2°. Пусть функционал F определен в банаховом пространстве B , обладает следующими свойствами: 1) оператор $P = \text{grad } F$ имеет в любой точке u производную Гато P'_u , удовлетворяющую неравенствам (УП.1.4-5); 2) $P(0) = 0$. Рассмотрим задачу о минимуме функционала $F(x) = (f, x)$, где f - фиксированный элемент

мент пространства \mathcal{B}^* . Пусть координатная система имеет вид (УП.1.6). Составим матрицу M_n чисел $[y_{nj}, y_{nk}]$; $j, k = 1, 2, \dots$, где кр. скобки означают скалярное произведение в B_0 - энергетическом пространстве оператора P_n (формула (УП.1.5)). Пусть $\lambda_n^{(0)}$ - наименьшее собственное число матрицы M_n , пусть $\gamma(n)$ - положительная функция натурального числа n , удовлетворяющая неравенству

$$\lambda_1^{(n)} > C\gamma^2(n); \quad C = \text{const.} \quad (\text{УП.3.1})$$

Введем N -мерные гильбертовы пространства \bar{X}_n, \bar{Y}_n с нормами

$$\forall b \in \bar{X}_n, \|b\|_{\bar{X}_n} = \gamma(n) \|b\|_{R_n}, \quad \|b\|_{\bar{Y}_n} = \frac{1}{\gamma(n)} \|b\|_{R_n}. \quad (\text{УП.3.2})$$

Примем, что координатная система удовлетворяет обычным условиям: 1) элементы $y_{n1}, y_{n2}, \dots, y_{nn}$ образуют базис в выделенном на \bar{X}_n подпространстве $B_0^{(n)}$ пространства B_0 ; 2) последовательность $\{B_0^{(n)}\}$ полна в B_0 . Добавим к этому еще одно условие более ограничительное, чем в линейном случае: 3) $y_{nk} \in \mathcal{D}(P)$. Как всегда в методе Рунге, будем искать приближенное решение задачи о минимуме как элемент подпространства $B_0^{(n)}$:

$$x^{(n)} = \sum_{k=1}^N a_k^{(n)} y_{nk}.$$

Условие 3) позволяет записать уравнения Рунге в форме Бубнова-Галеркина:

$$\left(P \left(\sum_{k=1}^N a_k^{(n)} y_{nk} \right), y_{nj} \right) = (f, y_{nj}), \quad 1 \leq j \leq N. \quad (\text{УП.3.3})$$

Обозначим $a^{(n)} = (a_1^{(n)}, a_2^{(n)}, \dots, a_N^{(n)})$, $f^{(n)} = ((f, y_{n1}), (f, y_{n2}), \dots, (f, y_{nN}))$, (f, y_{nr}) . Тогда эти

$$\left(P \left(\sum_{k=1}^N a_k^{(n)} y_{nk} \right), y_{nj} \right) = P_j^{(n)}(a^{(n)}),$$

$$P_n(a^{(n)}) = (P_1^{(n)}(a^{(n)}), P_2^{(n)}(a^{(n)}), \dots, P_N^{(n)}(a^{(n)})).$$

Вектор уравнений Рунге можно записать в виде

$$P_n(a^{(n)}) = f^{(n)}; \quad (\text{УП.3.4})$$

для рассуждать $a^{(n)}$ и $f^{(n)}$ как элементы пространства \bar{X}_n и

\bar{Y}_n соответственно.

Составление оператора P_n и вектора $f^{(n)}$ сопряжено с погрешностями вычислений. Пусть в результате вычислений мы вместо P_n и $f^{(n)}$ получили именные объекты $P_n + \Gamma_n$ и $f^{(n)} + \delta^{(n)}$. Тогда вместо системы (УП.3.4) мы на самом деле будем решать систему

$$(P_n + \Gamma_n)C^{(n)} = f^{(n)} + \delta^{(n)}. \quad (\text{УП.3.5})$$

3°. Теорема УП.3.1. В условиях данного параграфа процесс (УП.3.3) устойчив в (\bar{X}_n, \bar{Y}_n) .

Прежде всего покажем, что нормы $\|a^{(n)}\|_{\bar{X}_n}$ ограничены независимо от n . Так как

$$P(x_*^{(n)}) = P(x_*^{(n)}) - P(0) = \int_0^1 P'_{tx_*^{(n)}} x_*^{(n)} dt,$$

то уравнениям Рунге можно придать вид

$$\int_0^1 (P'_{tx_*^{(n)}} x_*^{(n)}, y_{nj}) dt = (f, y_{nj}), \quad 1 \leq j \leq N.$$

Умножая это на $a_j^{(n)}$ и складывая, найдем

$$\int_0^1 (P'_{tx_*^{(n)}} x_*^{(n)}, a_*^{(n)}) dt = (f, a_*^{(n)}).$$

Слева заменим подынтегральную функцию по неравенству (УП.1.5) меньшей величиной: $\alpha^2 (P'_{t_0} a_*^{(n)}, a_*^{(n)})$, а справа заменим $(f, a_*^{(n)})$ большей величиной $\|f\| \cdot \|a_*^{(n)}\|$:

$$\alpha^2 |a_*^{(n)}|^2 \leq \|f\| \cdot \|a_*^{(n)}\|,$$

т.е. норма в B_n . По неравенству (УП.1.4) $\|a_*^{(n)}\| \leq \gamma^{-1} \|a_*^{(n)}\|$, поэтому

$$|a_*^{(n)}| \leq \frac{\|f\|}{\alpha^2 \gamma}. \quad (\text{УП.3.6})$$

Итак,

$$\begin{aligned} |a_*^{(n)}|^2 &= (\text{Min}_{t \in I} (a^{(n)}, a^{(n)}))_{R_n} \geq \lambda_n \|a^{(n)}\|_{R_n}^2 \geq \\ &\geq C \gamma^2(n) \|a^{(n)}\|_{R_n}^2 = C \|a^{(n)}\|_{\bar{X}_n}^2 \end{aligned}$$

и, следовательно,

$$\|a^{(n)}\|_{\bar{x}_r} \leq \frac{1}{\sqrt{c}} |a_*^{(n)}| \leq \frac{\|P\|}{\alpha^2 \sqrt{c}}. \quad (\text{УЛ.3.7})$$

4°. Дальнейшее основано на результатах Лангенбаха [1], которые мы здесь коротко опишем. Пусть P оператор, вообще говоря нелинейный, определенный на линейном множестве $\mathcal{D}(P)$, плотном в банаховом пространстве B , и действующий из B в B^* . Пусть $P(0) = 0$. Предположим, что производная Гато P_u определена при любом $u \in \mathcal{D}(P)$ и $\mathcal{D}(P_u) \supset \mathcal{D}(P)$. Допустим далее, что эта производная удовлетворяет неравенствам (УЛ.1.4-5). Оператор P совпадает на $\mathcal{D}(P)$ с градиентом функционала

$$F(x) = \int_0^1 (P(tx), x) dt$$

и, следовательно, решение (если оно существует) уравнения

$$Px = f, \quad f \in B^* \quad (\text{УЛ.3.8})$$

реализует минимум функционала

$$F(x) - (f, x) = \int_0^1 (P(tx), x) dt - (f, x). \quad (\text{УЛ.3.9})$$

Эта последняя задача имеет обобщенное решение при любом $f \in B^*$ (см. § I данной главы); обозначим его через $x = \tilde{P}^{-1}f$. Нетрудно доказать, что различным x соответствуют различные f (см. Лангенбах I); существует, следовательно, обратный к \tilde{P}^{-1} оператор $(\tilde{P}^{-1})^{-1} = \tilde{P}$, и $\tilde{P}x = f$. Если $x \in \mathcal{D}(P)$, то $Px = f$. Это значит, что \tilde{P} есть расширение оператора P . Доказывается, что \tilde{P}^{-1} удовлетворяет условию Липшица

$$\|\tilde{P}^{-1}f_1 - \tilde{P}^{-1}f_2\| \leq \frac{1}{\alpha r} \|f_1 - f_2\|_{B^*}. \quad (\text{УЛ.3.10})$$

Рассмотрим теперь уравнение

$$\tilde{P}(x) + \Gamma(x) = f, \quad (\text{УЛ.3.11})$$

где Γ удовлетворяет условию

$$\|\Gamma(x_1) - \Gamma(x_2)\| \leq \alpha \|x_1 - x_2\|. \quad (\text{УЛ.3.12})$$

Доказывается, что при $\alpha < \alpha r$ уравнение (УЛ.3.10) имеет единственное решение; если x_* и \tilde{x}_* - решения уравнений (УЛ.3.

т (УЛ.3.11) соответственно, то $\|z_n - x_n\|_{R_n} \rightarrow 0$.

5°. Вернемся к теореме УЛ.3.1. Воспользуемся результатами, изложенными в п.4°. С этой целью отождествим \bar{P} с P_n , Γ - с Γ_n , B - с X_n ; тогда B^* следует отождествить с Y_n и B_0 - с R_n . Далее, отождествим x с $a^{(n)}$, f - с $f^{(n)}$ и z с $C^{(n)}$.

Оценим значения величин γ , α , α . Прежде всего,

$$P'_{n,a^{(n)}} b^{(n)} = (P'_{1,a^{(n)}} b^{(n)}, P'_{2,a^{(n)}} b^{(n)}, \dots, P'_{N,a^{(n)}} b^{(n)}).$$

Имеем

$$P'_{j,a^{(n)}} b^{(n)} = \frac{d}{dt} (P(x^{(n)} + ty^{(n)}), y_{nj})|_{t=0};$$

здесь

$$x^{(n)} = \sum_{k=1}^N a_k^{(n)} y_{nk}, \quad y^{(n)} = \sum_{k=1}^N b_k^{(n)} y_{nk}.$$

Далее

$$\begin{aligned} & \frac{d}{dt} (P(x^{(n)} + ty^{(n)}), y_{nj})|_{t=0} = \\ & = \left(\frac{d}{dt} P(x^{(n)} + ty^{(n)})|_{t=0}, y_{nj} \right) = (P'_{x^{(n)}} y^{(n)}, y_{nj}). \end{aligned}$$

Умножая это на $b_j^{(n)}$ и складывая, получим

$$(P'_{n,a^{(n)}} b^{(n)}, b^{(n)})_{R_n} = (P'_{x^{(n)}} y^{(n)}, y^{(n)}). \quad (\text{УЛ.3.13})$$

По неравенству (УЛ.1.5)

$$(P'_{n,a^{(n)}} b^{(n)}, b^{(n)})_{R_n} \geq \alpha^2 (P'_{y_0} y^{(n)}, y^{(n)}) = \alpha^2 \|y^{(n)}\|^2. \quad (\text{УЛ.3.14})$$

Воспользуемся теперь неравенством (УЛ.3.7):

$$(P'_{n,a^{(n)}} b^{(n)}, b^{(n)})_{R_n} \geq C \alpha^2 \|b^{(n)}\|_{X_n}. \quad (\text{УЛ.3.15})$$

Неравенство (УЛ.3.14) является аналогом неравенства (УЛ.1.4): отсюда видно, что в утверждение п.4° можно, применительно к нашему случаю, заменить γ на $\alpha \sqrt{C}$; важно отметить, что последняя величина не зависит от n . Далее, по формуле (УЛ.3.15) -14)

$$\begin{aligned} (P'_{n,a^{(n)}} b^{(n)}, b^{(n)})_{R_n} &= (P'_{x^{(n)}} y^{(n)}, y^{(n)}) \geq \\ &\geq d^2 (P'_0 y^{(n)}, y^{(n)}) = d^2 (P'_{n,0} b^{(n)}, b^{(n)}). \end{aligned}$$

(УП.3.16)

Мы приняли здесь $U_0 = 0$ что, очевидно, не уменьшает общности. Неравенство (УП.3.16) есть аналог неравенства (УП.1.5) с той же постоянной d .

Из результатов, перечисленных в п.4⁰, вытекает теперь, что операторы $P'_{n,a^{(n)}}$ удовлетворяют условию Липшица с постоянной C , которая не зависит от n и от $f^{(n)}$.

Рассмотрим теперь уравнение Гитца (УП.3.4), искаженное уравнение Рунца (УП.3.5) и вспомогательное уравнение

$$P_n q^{(n)} = f^{(n)} + \delta^{(n)}. \quad (\text{УП.3.17})$$

Так как P'_n удовлетворяет условию Липшица с постоянной C ,

$$\|q^{(n)} - a^{(n)}\|_{X_n} \leq C \|\delta^{(n)}\|_{Y_n}, \quad \delta \rightarrow 0. \quad (\text{УП.3.18})$$

Далее, P_n удовлетворяет условию Липшица (УП.3.12), в котором следует положить $\alpha = \Psi(t, \delta)$. Как было доказано в п.3⁰, но мы $\|a^{(n)}\|_{X_n}$ ограничены, поэтому можно считать число t фиксированным; в силу сказанного в конце п.4⁰, $\|C^{(n)} - q^{(n)}\|_{X_n} \xrightarrow{\delta \rightarrow 0} 0$. Вместе с неравенством (УП.3.18) это дает

$$\|C^{(n)} - a^{(n)}\|_{X_n} \xrightarrow{\delta \rightarrow 0} 0. \quad (\text{УП.3.19})$$

а это означает, что вычислительный процесс (УП.3.3) устойчив в последовательности (X_n, Y_n) .

Замечание УП.3.1. Если координатная система сильно минимальна в энергетической норме оператора P'_0 , то можно положить $f(n) = I$. Отсюда следует, что в рассматриваемом случае процесс (УП.3.3) устойчив в последовательности (R_n, R_n) .

6⁰. Теорема УП.3.2. Если координатная система почти ортонормирована в B_0 , то процесс вычисления приближенного решения по Рунцу условия настоящего параграфа устойчив в последова-

тельности $(B_0^{(n)}, R_{X'})$.

Доказательство очень просто. Если $x_*^{(n)}$ - искомого приближенного решение по Рунту, то

$$\begin{aligned} |x_*^{(n)} - x_*^{(n)}| &= (M_n(c^{(n)} - a^{(n)}), c^{(n)} - a^{(n)})_{R_{X'}} \leq \\ &\leq \Lambda_0 \|c^{(n)} - a^{(n)}\|_{R_{X'}}^2, \end{aligned}$$

где Λ_0 - верхняя граница собственных чисел матрицы M_n . Почти ортонормированная система сильно минимальна, и процесс гильбертовых коэфф. Фурье Рунта устойчив в $(R_{X'}, R_{X'})$, и при достаточно малом δ будет $\|c^{(n)} - a^{(n)}\|_{R_{X'}} \leq \varepsilon$. Но тогда $|x_*^{(n)} - x_*^{(n)}| \leq \sqrt{\Lambda_0} \delta$, и теорема доказана.

7°. Обратимся к МКЭ: для определенности остановимся на варианте этого метода, развитом в книге автора [22]. В этом случае система координатных функций не сильно минимальна в энергетической норме: в частности, в случае первой краевой задачи, а также в случае кубической области и более произвольных краевых условий, наименьшее собственное число матрицы Грама элементов (УИ.1.6) имеет порядок $O(h^m)$, где h - шаг сетки и m - размерность евклидова пространства задачи (см. [22], гл.11, §§ 2,3¹⁾).

Справедлива теорема, которая является частным случаем теоремы УИ.3.1:

Теорема УИ.3.3. Пусть достаточно гладкий функционал $F(x)$ определен на функциях пространства $V = W_p^{(s)}(\Omega)$, где Ω - ограниченная область в R_m . Пусть операторы $P = \text{grad} F$ и P'_u определены на множестве, плотном в V , причем P'_u удовлетворяет неравенствам (УИ.1.4-5). Пусть $h = 1/2n$ - шаг сетки, и $N = N(n)$ - число координатных функций, входящих в выражение приближенного решения. Вычислительный процесс МКЭ устойчив в последовательности пар N -мерных гильбертовых пространств (X_n, Y_n) с норма-

$$\forall a \in R_n, \|a^{(n)}\|_{X_n} = h^{m/2} \|a^{(n)}\|_{R_n}, \|a^{(n)}\|_{Y_n} = h^{-m/2} \|a^{(n)}\|_{R_n}.$$

Аналогичная теорема верна и для задачи с естественными краевыми условиями, если Ω - есть куб в R_m .

1) по надобности автора в с. 165 вместо начального h^{-m} должно h^m .

§ 4. Точности искажения и округления рекуррентного вычислительного процесса

1°. Рассмотрим нелинейный рекуррентный процесс

где $x^{(0)} \in X_n$, $f^{(n)} \in Y_n$, а X_n, Y_n - банаховы пространства. Искраженный процесс запишем в виде

$$A_n(x^{(0)}, x^{(1)}, \dots, x^{(n)}) = f^{(n)}; \quad n = 0, 1, 2, \dots, \quad (\text{УП.4.1})$$

$$A_n(\bar{x}^{(0)}, \bar{x}^{(1)}, \dots, \bar{x}^{(n)}) + \Gamma_n(\bar{x}^{(0)}, \bar{x}^{(1)}, \dots, \bar{x}^{(n)}) = f^{(n)} + \delta^{(n)}; \quad n = 0, 1, 2, \dots \quad (\text{УП.4.2})$$

Примем следующие допущения: 1) если элементы $x^{(0)}, x^{(1)}, \dots, x^{(n)}$ фиксированы и их нормы не превосходят некоторого числа \bar{t} , то n -е уравнение (УП.4.1) имеет не более одного решения в шаре $\|x^{(n)}\| \leq \bar{t}$; 2) существует такое число $\bar{T} = \text{const}$, что если $\|f^{(n)}\| \leq \bar{T}$ при любом n , то бесконечная система (УП.4.1) разрешима и $\|x^{(n)}\|$

нетрудно указать свободный вычислительный процесс, эквивалентный рекуррентному процессу (УП.4.1). Построим две новые пространства \bar{X}_n и \bar{Y}_n . $n = 0, 1, 2, \dots$ Элементами пространства \bar{X}_n являются конечные последовательности $\bar{x}^{(n)} = (x^{(0)}, x^{(1)}, \dots, x^{(n)})$, $x^{(k)} \in X_k$; норму в \bar{X}_n определим формулой

$$\|\bar{x}^{(n)}\|_{\bar{X}_n} = \max_{k \leq n} \|x^{(k)}\|_{X_k} \quad (\text{УП.4.3})$$

(допустимы и другие определения нормы); аналогично определим \bar{Y}_n . Определим еще оператор \bar{A}_n , действующий из \bar{X}_n в \bar{Y}_n по формуле

$$\bar{A}_n(\bar{x}^{(n)}) = \{A_k(x^{(0)}, \dots, x^{(k)})\}_{k=0}^n; \quad (\text{УП.4.4})$$

аналогично определим оператор $\bar{\Gamma}_n$. Уравнения (УП.4.1) и (УП.4.2) соответственно равносильны уравнениям свободных вычислительных процессов

$$\forall n \geq 0, \bar{A}_n(\bar{x}^{(n)}) = \bar{f}^{(n)}; \quad (\text{УП.4.5})$$

$$\forall n \geq 0, \bar{A}_n(\bar{x}^{(n)}) + \bar{\Gamma}_n(\bar{x}^{(n)}) = \bar{f}^{(n)} + \bar{\delta}^{(n)}. \quad (\text{УП.4.6})$$

Мы дадим здесь определение устойчивости нелинейного рекуррентного процесса, отличное от аналогичного определения для линейного процесса, данного в гл. III. Рекуррентный процесс (УП.4.1) назовем устойчивым, если устойчив соответствующий свободный процесс (УП.4.5).

2°. Условия устойчивости рекуррентного процесса получаются из результатов § 2. Эти условия таковы: 1) $A_n(0) = 0$; 2) оператор \bar{A}_n^{-1} определен в шаре $\|f^{(n)}\| \leq T$ и удовлетворяет условию Липшица $\|\bar{A}_n^{-1}(f_1) - \bar{A}_n^{-1}(f_2)\| \leq C \|f_1 - f_2\|$

с постоянной, не зависящей от n ; 3) $\|\delta^{(n)}\| < \delta$; 4) если $\bar{x}, \bar{x}', \bar{x}'' \in X_n$ и нормы этих элементов не превосходят t , то $\|\bar{\Gamma}_n(\bar{x})\| < \psi(t, \delta)$; $\|\bar{\Gamma}_n(\bar{x}) - \bar{\Gamma}_n(\bar{x}')\| < \psi(t, \delta) \|\bar{x}' - \bar{x}\|$.

Переформулируем эти условия в терминах операторов A_n и Γ_n . Условие 1) означает, что

$$\forall n \geq 0, A_n(0, 0, \dots, 0) = 0. \quad (\text{УП.4.7})$$

Далее, \bar{A}_n^{-1} есть оператор, определяющий решение конечной системы

$$A_k(x^{(0)}, x^{(1)}, \dots, x^{(k)}) = f^{(k)}, \quad k = 0, 1, \dots, n.$$

Это решение имеет вид

$$\bar{x}^{(k)} = B_k(f^{(0)}, f^{(1)}, \dots, f^{(k)}), \quad k = 0, 1, \dots, n.$$

Отсюда следует, что $\bar{A}_n^{-1} = \{B_0, B_1, \dots, B_n\}$, и условие 2) принимает вид

$$\forall n \geq 0, \|B_n(f_1^{(0)}, f_1^{(1)}, \dots, f_1^{(n)}) - B_n(f_2^{(0)}, f_2^{(1)}, \dots, f_2^{(n)})\| <$$

$$< C \max_k \|f_1^{(k)} - f_2^{(k)}\|.$$

(УП.4.8)

Условие 3) означает, что

$$\|\delta^{(n)}\| < \delta.$$

(УП.4.9)

Наконец, условия 4) своятся к следующему:

$$\|\Gamma_n(x^{(0)}, x^{(1)}, \dots, x^{(n)})\| < \psi(t, \delta)$$

(УП.4.10а)

$$\| \Gamma_n(x_1^{(0)}, x_1^{(1)}, \dots, x_1^{(n)}) - \Gamma_n(x_2^{(0)}, x_2^{(1)}, \dots, x_2^{(n)}) \| \leq \leq \Psi(t, \delta) \max_{0 \leq k \leq n} \| x_1^{(k)} - x_2^{(k)} \|. \quad (\text{УП.4.1})$$

На основании теоремы УП.2.1 можно утверждать, что найденный рекуррентный процесс (УП.4.1) устойчив в смысле определения настоящего параграфа, если $\| f^{(n)} \| \leq T = \text{const}$ и выполнены условия (УП.4.7-10).

Если процесс (УП.4.1) устойчив, то при малых δ погрешность нахождения ограничена независимо от n .

3^c. Укажем случай, когда погрешность округления рекуррентного валиейного процесса также остается ограниченной. Обозначим $A_n + \Gamma_n = \tilde{A}_n$, $f^{(n)} + \delta^{(n)} = \tilde{f}^{(n)}$, и пусть точное решение искомого уравнения (УП.4.2) имеет вид

$$\forall n \geq 0, \tilde{x} = \tilde{B}_n(\tilde{x}^{(0)}, \tilde{x}^{(1)}, \dots, \tilde{x}^{(n-1)}, \tilde{f}^{(n)}). \quad (\text{УП.4.1})$$

допустим, что операторы \tilde{B}_n удовлетворяют условию Липшица

$$\| \tilde{B}_n(\tilde{x}_1^{(0)}, \tilde{x}_1^{(1)}, \dots, \tilde{x}_1^{(n-1)}, \tilde{f}^{(n)}) - \tilde{B}_n(\tilde{x}_2^{(0)}, \tilde{x}_2^{(1)}, \dots, \tilde{x}_2^{(n-1)}, \tilde{f}^{(n)}) \| \leq \leq q \max_{0 \leq k \leq n} \| \tilde{x}_1^{(k)} - \tilde{x}_2^{(k)} \|, \quad q = \text{const} < 1. \quad (\text{УП.4.1})$$

Пусть $v^{(0)}, v^{(1)}, \dots, v^{(n-1)}$ - точно известные элементы, $\tilde{x}^{(n)}$ - точное значение выражения $\tilde{B}_n(v^{(0)}, v^{(1)}, \dots, v^{(n-1)})$.

Допустим, что в результате погрешностей вычисления мы вместо $\tilde{x}^{(n)}$ получим элемент $u^{(n)} = \tilde{x}^{(n)} + d^{(n)}$. Будем считать, что

погрешность округления ограничена числом $\varepsilon > 0$, так что $\| d^{(n)} \| \leq \varepsilon \| v^{(n)} \|$. По неравенству (УП.4.12) имеем

$$\forall n \geq 0, \| u^{(n)} - \tilde{x}^{(n)} \| \leq q \max_{0 \leq k \leq n-1} \| v^{(k)} - \tilde{x}^{(k)} \| + \varepsilon \| v^{(n)} \|.$$

Допустим, что $\| \tilde{x}^{(0)} \| \leq \text{const}$, получаем из последнего неравенства

$$\hat{\tilde{x}}_n := \| v^{(n)} - \tilde{x}^{(n)} \| \leq q \max_{0 \leq k \leq n-1} \hat{\tilde{x}}_k + C\varepsilon, \quad q = \frac{q}{1-\varepsilon}. \quad (\text{УП.4.12})$$

Пусть $\max_{0 \leq k \leq n} \hat{\xi}_k^{(n)} = \hat{\xi}^{(n)}$, $0 \leq l \leq n$. Неравенство (УЛ.4.13) дает при $n = l$ неравенство $\hat{\xi}_l \leq q' \hat{\xi}_l + C\varepsilon$ или, если $q' < 1$, $\hat{\xi}_l \leq C\varepsilon / (1 - q')$. Отсюда следует искомая оценка погрешности округления

$$\hat{\xi}_n \leq \frac{C\varepsilon}{1 - q'}. \quad (\text{УЛ.4.14})$$

§ 5. Метод Ньютона - Канторовича

10. В настоящем параграфе мы исследуем влияние искажений и округления двух методов Ньютона - Канторовича: модифицированного и основного. Эти методы приводят к рекуррентным процессам более специального вида, нежели рассмотренные в § 4. Устойчивость таких процессов можно исследовать проще и в более общих предположениях. С рассмотрения процессов такого вида мы и начнем.

Пусть задан рекуррентный процесс вида

$$A_n x^{(n)} + Q_n(x^{(n-1)}) = f^{(n)}, \quad (\text{УЛ.5.1})$$

где A_n - линейные, а Q_n - нелинейные операторы. Для простоты допустим, что все эти операторы действуют в одном и том же банаховом пространстве X . Искаженный процесс запишем в виде

$$A_n \tilde{x}^{(n)} + \Gamma_n \tilde{x}^{(n)} + Q(\tilde{x}^{(n-1)}) + \Gamma_{2n}(\tilde{x}^{(n-1)}) = f^{(n)} + \delta^{(n)}. \quad (\text{УЛ.5.2})$$

Примем следующие допущения. 1) Операторы A_n^{-1} равномерно ограничены: $\|A_n^{-1}\| \leq a = \text{const}$. 2) При любом n оператор $A_n^{-1} Q_n$ есть оператор сжатия, так что

$$\|A_n^{-1} Q_n(x) - A_n^{-1} Q_n(y)\| \leq q \|x - y\|, \quad q = \text{const} < 1. \quad (\text{УЛ.5.3})$$

3) Справедливы следующие неравенства:

$$\|\delta^{(n)}\| \leq \delta; \quad (\text{УЛ.5.4})$$

если $\|x\|, \|x'\|, \|x''\| \leq t$, то (к 1, 2)

$$\|A_n^{-1} \Gamma_{kn}(x)\| \leq \psi(t, \delta), \quad (\text{УЛ.5.5})$$

$$\|A_n^{-1} \Gamma_{kn}(x) - A_n^{-1} \Gamma_{kn}(x')\| \leq \psi(t, \delta) \|x - x'\|; \quad (\text{УЛ.5.6})$$

Функции φ и ψ обладают теми же свойствами, что и в предыдущих параграфах.

Имеем

$$x^{(n)} = A_n^{-1} f^{(n)} - A_n^{-1} Q_n(x^{(n-1)}),$$

$$\tilde{x}^{(n)} = A_n^{-1} f^{(n)} + A_n^{-1} \delta^{(n)} - A_n^{-1} Q_n(\tilde{x}^{(n-1)}) -$$

$$- A_n^{-1} \Gamma_{1n} \tilde{x}^{(n)} - A_n^{-1} \Gamma_{2n}(\tilde{x}^{(n)}).$$

Отсюда

$$\tilde{x}^{(n)} - x^{(n)} =$$

$$= A_n^{-1} \delta^{(n)} - [A_n^{-1} Q_n(\tilde{x}^{(n-1)}) - A_n^{-1} Q_n(x^{(n-1)})] - A_n^{-1} \Gamma_{1n}(\tilde{x}^{(n)}) - A_n^{-1} \Gamma_{2n}(\tilde{x}^{(n-1)}).$$

Как и в предыдущих параграфах, доказываем, что если нормы $\|x^{(n)}\|$ ограничены в совокупности, $\|x^{(n)}\| \leq t$, то уравнение (УЛ.5.2) имеет в шаре $\|\tilde{x}\| \leq 2t$ одно и только одно решение. Оценим разность $\tilde{x}^{(n)} - x^{(n)}$, где $x^{(n)}$ и $\tilde{x}^{(n)}$ - решения уравнений (УЛ.5.1) и (УЛ.5.2) соответственно. Имеем

$$\|\tilde{x}^{(n)} - x^{(n)}\| \leq q \|\tilde{x}^{(n-1)} - x^{(n-1)}\| + a\delta + 2\varphi(2t, \delta),$$

или, обозначая $\|\tilde{x}^{(n)} - x^{(n)}\| = \hat{\xi}_n$, $a\delta + 2\varphi(2t, \delta) = \eta$,

$$\hat{\xi}_n \leq q \hat{\xi}_{n-1} + \eta.$$

Известно, что $\hat{\xi}_n \leq \tau_n$, где τ_n - решение системы $\tau_0 = \hat{\xi}_0$, $\tau_n = q\tau_{n-1} + \eta$, $n = 1, 2, \dots$. Это решение есть $\tau_n = q^n \hat{\xi}_0 + \eta(1-q)^{-1}$ и, следовательно

$$\|\tilde{x}^{(n)} - x^{(n)}\| \leq q^n \|\tilde{x}^{(0)} - x^{(0)}\| + (1-q)^{-1} [a\delta + 2\varphi(2t, \delta)]. \quad (\text{УЛ.5.7})$$

Неравенство (УЛ.5.7) дает оценку погрешности искажения для процесса (УЛ.5.1); из этого неравенства вытекает также устойчи-

ность этого процесса при условиях (УП.5.3-6).

Замечание УП.5.1. Джиккарони [8, Ю, II] рассмотрел частный случай процесса (УП.5.1), который получается применением классического метода Рунда к уравнению $Ax + Qx = f$, где A - самосопряженный положительно определенный оператор с дискретным спектром, действующий в сепарабельном гильбертовом пространстве H , а Q - нелинейный непрерывный потенциальный оператор, имеющий положительно определенный дифференциал Фреше. Предполагается, что Q определен на всем пространстве H и $Q(0) = 0$. Наконец, $f \in H$. Принимаем, что рассматриваемое уравнение имеет единственное решение.

В качестве координатных элементов выбираются собственные элементы оператора B , сходного (см. книгу автора [8]) с A . Доказывается, что невяка приближения ε решения по Рунду, $x_*^{(n)}$, стремится к нулю:

$$\varepsilon^{(n)} = \|Ax_*^{(n)} - Qx_*^{(n)} - f\|_H \rightarrow 0.$$

Доказывается, что в метрике H погрешность аппроксимации имеет оценку (ω_{n+1} есть $(n+1)$ -е собственное число оператора B)

$$\|x_* - x_*^{(n)}\|_H = O\left(\frac{\|\delta^{(n)}\|_B}{\omega_{n+1}}\right),$$

а в энергетической норме оператор A - оценку

$$|x_* - x_*^{(n)}| = O\left(\frac{\|\delta^{(n)}\|}{\sqrt{\omega_{n+1}}}\right).$$

Устойчивость процесса доказана в том же предположении, что координатные элементы суть собственные элементы оператора, сходного с A .

2°. Переходим к ошибкам округления. Обозначим $A_n = \Gamma_{1n} = \tilde{A}_n$, $Q_n + \Gamma_{2n} = \tilde{Q}_n$ - оба эти оператора точно известны, и оператор \tilde{A}_n существует. Допустим, что оператор $\tilde{A}_n \tilde{Q}_n$ есть оператор скалярного умножения на некотором шаре, содержащем $x^{(n)}$, с коэффициентом скалярного умножения q_n , зависящим от n . Уравнение (УП.5.2) можно привести вид

$$x^{(n)} = a_n (x^{(n-1)})^2 + b^{(n)},$$

$$a_n = \tilde{A}_n^{-1} \tilde{Q}_n, \quad b^{(n)} = \tilde{A}_n^{-1} (C_n^2 + \delta^{(n)}). \quad (2.1.1)$$

Уравнение (УЛ.5.8) разрешимо итерациями, которые можно записать так:

$$m=0,1,2,\dots; y^{(m)} = a_n(y^{(m-1)}) + b^{(m)}, \quad (\text{УЛ.5.9})$$

где $y^{(m)}$ есть m -е приближение к решению $x^{(n)}$. Допустим, что, приняв правую часть формулы (УЛ.5.9), мы совершаем ошибку округления, относительная величина которой не превосходит числа $\varepsilon > 0$. Пусть, далее, известное нам округленное значение $(m-1)$ -го приближения $y^{(m-1)}$ есть $v^{(m-1)}$. Тогда в результате вычислений мы вместо элемента $y^{(m)}$ получим некоторый элемент $v^{(m)}$, удовлетворяющий соотношению

$$v^{(m)} = a_n(v^{(m-1)}) + b^{(m)} + d^{(m)}, \quad (\text{УЛ.5.10})$$

$$\|d^{(m)}\| \leq \varepsilon \|v^{(m)}\|.$$

Вычитая отсюда равенство (УЛ.5.9), получаем

$$v^{(m)} - y^{(m)} = a_n(v^{(m-1)}) - a_n(y^{(m-1)}) + d^{(m)}$$

и

$$\gamma_m := \|v^{(m)} - y^{(m)}\| \leq \tilde{q}_n \gamma_{m-1} + \varepsilon \gamma_m + \varepsilon \|y^{(m-1)}\|.$$

Допустим, что нормы $\|b^{(m)}\|$ ограничены в совокупности: $\|b^{(m)}\| \leq C$. Тогда нормы $\|y^{(m)}\|$ также ограничены; пусть $\|y^{(m-1)}\| \leq C_0$. В этом случае

$$\gamma_m \leq \frac{\tilde{q}_n}{1-\varepsilon} \gamma_{m-1} + \frac{C_0 \varepsilon}{1-\varepsilon}. \quad (\text{УЛ.5.11})$$

Обозначим $\tilde{q}_n/(1-\varepsilon) = \nu$, $C_0 \varepsilon/(1-\varepsilon) = \lambda$, так что $\gamma_m \leq \nu \gamma_{m-1} + \lambda$. Заметим, что $v^{(0)} = y^{(0)}$; так как $y^{(0)}$ — начальное значение — не вычислялось, а задается, $\gamma_0 = 0$. Теперь из (УЛ.5.11) следует

$$\gamma_m \leq \frac{1-\nu^m}{1-\nu} \frac{C_0 \varepsilon}{1-\varepsilon}. \quad (\text{УЛ.5.12})$$

Если ε столь мало, что $\nu < 1$, то оценка (УЛ.5.12) дает далее

$$\gamma_m \leq \frac{C_0 \varepsilon}{(1-\nu)(1-\varepsilon)}; \quad (\text{УЛ.5.13})$$

ошибка округления ограничена независимо от m и n и равномерно

стремится по этим индукциям к нулю при $\varepsilon \rightarrow 0$. Если ε не очень мало, так что $\nu > 1$, то оценка (УП.5.12) стремится к бесконечности вместе с m .

3°. Модифицированный метод Ньютона - Канторовича применяется к решению нелинейных задач вида $P(x) = 0$. Достаточно полное описание этого метода, условия и скорость его сходимости даны в книге Канторовича и Акилова [1]; здесь мы отметим только, что сам метод состоит в применении рекуррентного процесса

$$x^{(n)} = x^{(n-1)} - A^{-1} P(x^{(n-1)}); \quad A = P'_{x^{(0)}}; \quad (\text{УП.5.14})$$

$x^{(0)}$ - начальное приближение к искомому решению. Нам будет удобно записать этот процесс в виде

$$A x^{(n)} + Q(x^{(n-1)}) = 0; \quad Q = A - P. \quad (\text{УП.5.15})$$

Заметим, что оператор A не задан, а вычисляется, естественно, с некоторой погрешностью.

Проверим выполнение условий п.1°. Условие $\|A_n^{-1}\| \leq \text{const}$ вытекает из того, что оператор $A_n = A = P'_{x^{(0)}}$ не зависит от n , а что ограниченность оператора $[P'_{x^{(0)}}]$ принадлежит к числу условий сходимости модифицированного метода Ньютона - Канторовича. Условие (УП.5.3) будет выполнено, если P имеет непрерывную первую производную; в теореме о сходимости метода требуется существование второй производной. Мы предполагаем также, что выполнены условия (УП.5.4-6), характеризующие искажение процесса. Это показано в пп.1°, 2°, вычислительный процесс, связанный с модифицированным методом Ньютона - Канторовича, устойчив; если уравнения модифицированного процесса решать и решать с достаточно малыми относительными погрешностями, то погрешность решения остается ограниченной.

4°. В книге Канторовича и Акилова [1] доказано следующее утверждение. Пусть $x^{(0)}$ - начальное приближение, а $x^{(1)}, x^{(2)}, \dots, x^{(n)}, \dots$ - приближения, построенные по основному методу Ньютона - Канторовича. Если в точке $x^{(0)}$ выполнены условия сходимости модифицированного метода Ньютона - Канторовича, то они выполнены и в точках $x^{(1)}, x^{(2)}, \dots$. Отсюда нетрудно получить, что для вычислительного процесса, связанного с основным методом Ньютона - Канторовича, справедливы заключения, сделанные в конце п.3°. Для модифицированного метода.

ГЛАВА УШ

Некоторые односторонние вариационные задачи.

В обширной литературе по задачам, связанным с вариационными неравенствами, видное место занимают так называемые односторонние вариационные задачи. Обычная постановка этих задач такова: на рефлексивном банаховом пространстве B заданы полунепрерывный снизу возрастающий функционал $J(x)$ и замкнутое выпуклое множество M . Ставится задача: найти такой элемент $x_* \in M$, чтобы

$$\forall x \in M, J(x_*) \leq J(x).$$

Условие $x \in M$ называется односторонним ограничением, а множество M называется множеством ограничения.

В книге Дювоа [1] доказано, что такая задача имеет решение одно единственно, если функционал J - строго выпуклый. В книге Дювоа и Лионса [1] рассмотрен ряд задач физики и механики, которые приводятся к вариационным неравенствам и, в частности, к односторонним вариационным задачам. В книге Гловинского, Лисне и Тремольера [1] излагается ряд методов приближенного решения таких задач и даны приложения этих методов к большому числу более конкретных задач. Для многих приближенных методов даны доказательства сходимости. Однако, по мнению авторов этой книги, получение оценок погрешности приближенных методов - задача более деликатная.

В настоящей главе будут изложены результаты работ автора [26-30] по оценке погрешностей для односторонних вариационных задач; в основном речь будет идти о некотором классе таких задач, подробное описание которого дано ниже, в § 1. Будут изложены также (в § 4), но без доказательств, некоторые результаты работ Демьяновича [4] и Тухтына [1]. В конце главы будет дан обзор работ последнего времени по оценке погрешности аппроксимации для вариационных неравенств; этот обзор написан Тухтыным.

§ 1. Постановка задачи и ее приближенное решение.

1°. Пусть требуется решить задачу о минимуме функционала

$$J(x) - \lambda x \tag{1.1.1}$$

при одностороннем ограничении

$$\|x - q\| \leq \delta, \quad \delta = \text{const.}$$

где \mathcal{F} - квадратичная форма, определенная и положительная на некотором линейном множестве. Прежде чем переходить к описанию остальных данных задачи, заметим, что на упомянутом множестве можно построить энергетическое пространство H со скалярным произведением $[u, v] = \mathcal{F}(u, v)$ и нормой $\|u\| = [\mathcal{F}(u)]^{1/2}$; здесь $\mathcal{F}(u, v)$ - билинейная форма, соответствующая квадратичной форме $\mathcal{F}(u)$. Примем, что ℓ - линейный функционал, ограниченный в энергетической норме, и что $\|\cdot\|$ - полунорма, определенная (но не обязательно всюду конечная) на всем пространстве H .

Введем еще в рассмотрение гильбертово пространство \tilde{H} , по отношению к которому H является подпространством (может быть, несобственным). Скалярное произведение и норму в H будем обозначать теми же смыслами, что и в \tilde{H} . Будем считать, что g - заданный элемент пространства \tilde{H} и что полунорма $\|\cdot\|$ определена (хотя, разумеется, не обязательно всюду конечна) на всех элементах этого пространства. Необходимо потребовать, чтобы множество (УШ. I. 2) не было пустым - для этого необходимо и достаточно, чтобы $\|g\| < \infty$ и чтобы расстояние $\rho(g, H)$, вычисленное в метрике $\|\cdot\|$, превосходило числа d .

2°. Последний вопрос рассмотрим подробнее. Вычислим величину $\alpha = \rho(g, H)$. Если $\alpha > d$, то, как только-что было отмечено, множество

$$M = \{x \in H : \|x - g\| \leq d\} \quad (\text{УШ. I. 3})$$

пусто, и задача (УШ. I. 1-2) не имеет решения. Пусть $\alpha < d$, тогда существует такой элемент $x' \in H$, что $\|x' - g\| < d$, и множество M не пусто. Положим $g - x' = h$, $x - x' = y$, тогда множество ограниченный определится соотношением $\|y - h\| \leq d$, в котором $\|h\| < d$.

Пусть теперь $\alpha = d$. Может случиться, что $\inf_{x \in H} \|x - g\|$ не достигается, тогда $\forall x \in H, \|x - g\| > d$, и множество M пусто. Если же инфимум достигается, то множество M не пусто, и задача (УШ. I. 1-2) имеет решение (см. ниже, п. 3°). Однако, это решение устойчиво относительно малых изменений d . Действительно, если заменить d на $d' < d$, то $\alpha = \alpha > d'$, то же множество

$$M' = \{x \in H, \|x - g\| < d'\}$$

пусто, и новая односторонняя вариационная задача не имеет решения.

Итак, рассматривая множество вида (VII.1.3), мы всегда считаем, что $\beta := \inf \|g\| < d$.

3°. По теореме Рунда существует единственный элемент x_0 , удовлетворяющий тождеству $\forall x \in H, \langle x, x_0 \rangle = [x, g]$; этот элемент решающий задачу о минимуме функционала, (VII.1.1) при отсутствии ограничений, будем считать известным.

Задача (VII.1.1-2) равносильна следующей:

$$\|x - x_0\| = \min, \|x - g\| \leq d. \quad (\text{VII.1.1})$$

Множество ограничений M , определяемое неравенством (VII.1.3) (или (VII.1.3)) - выпуклое; мы принимаем полунорму $\|\cdot\|$ так, что это множество замкнуто в метрике H . В силу результатов изложенных в книге Дюнса [1], гл.1, задача (VII.1.4) имеет и только одно решение, которое мы обозначим через x_* . Если $\|x_0 - g\| \leq d$, то, очевидно, $x_* = x_0$; если же $\|x_0 - g\| > d$ x_* решает задачу

$$\|x_* - x_0\| = \min, \|x_* - g\| = d. \quad (\text{VII.1.1})$$

Докажем это. Рассмотрим прямую линию отрезок Λ , соединяющую точки x_0 и x_* : $\Lambda = \{x \in H: x = \lambda x_* + (1-\lambda)x_0; 0 \leq \lambda \leq 1\}$. Установим Λ направлением от x_* к x_0 и в соответствии с этим назовем x_* левым, а x_0 - правым концом отрезка Λ . Ясно, что левый конец отрезка Λ принадлежит множеству M , правый - множеству $H \setminus M$. Идем на Λ расстояния d , порожденное нормой $\|\cdot\|$, и разделим отрезок пополам. Один из двух отрезков деления таков, что его левый конец принадлежит M , а правый принадлежит $H \setminus M$. Этот отрезок разделим пополам, и т.д. В результате получим последовательность вложенных отрезков одинаковой длины, стремящихся к x_* , и обладающих тем свойством, что их левые концы принадлежат M , а правые - $H \setminus M$. Эти отрезки содержат ровно одну общую точку: обозначим ее через y . Она - предельная для множества так как M замкнуто, то $y \in M$, и $\|y - g\| \leq d$. Одновременно то же предельная для $H \setminus M$: если $\{y_k\}$ - последовательность правых концов рассмотренных нами вложенных отрезков, то $\|y_k - y\| \rightarrow 0$. Однако если отрезке полунорма $\|\cdot\|$ непрерывна, поэтому $\|y - y_k\| \rightarrow 0$ и $\|g - y_k\| > d$, следовательно, $\|g - y\| \geq d$. Следовательно

$$\|g - y\| = d.$$

Допустим теперь, что $\|x_* - g\| < d$. Тогда на отрезке Λ точна y - промежуточная между x_* и x_0 , поэтому $|x_0 - y| < |x_0 - x_*| = \min_{x \in M} |x_0 - x|$, а это нелепо. Отсюда следует, что $\|x_* - g\| = d$, и наше утверждение доказано.

Множество $\{x \in N : \|x - g\| = d\}$ назовем границей множества M и обозначим через ∂M . Можно, следовательно, считать, что при $\|x_0 - g\| > d$ элемент x_* решает задачу

$$|x - x_0| = \min, \quad x \in M. \quad (\text{УШ. I. 6})$$

Ниже ввиду принимается, что $\|x_0 - g\| > d$.

4°. Задачу (УШ. I. 6) будем решать приближенно, сводя ее к конечномерной задаче. Выберем натуральное число n и N -мерное ($N = N(n)$) подпространство H_n пространства H ; потребуем, чтобы на элементах подпространства H_n полунорма $\|\cdot\|$ была конечной. Поставим одностороннюю вариационную задачу

$$x^{(n)} \in H_n \cap M, \quad |x_0 - x^{(n)}| = \min. \quad (\text{УШ. I. 7})$$

Задачу (УШ. I. 7) можно несколько упростить. Пусть \mathcal{O}_{on} - ортогональная проекция элемента x_0 на H_n , тогда

$$|x_0 - x^{(n)}|^2 = |x_{on} - x^{(n)}|^2 + |x_0 - x_{on}|^2.$$

второе слагаемое справа - постоянное, поэтому задача (УШ. I. 7) равносильна следующей:

$$|x_{on} - x^{(n)}|^2 = \min, \quad x^{(n)} \in H_n \cap M. \quad (\text{УШ. I. 8})$$

Преимуществом этой новой формы нашей задачи является то, что искомый элемент x_{on} , так и искомый элемент $x^{(n)}$ принадлежат одному и тому же пространству H_n .

Пересечение подпространства H_n и замкнутого выпуклого множества M есть замкнутое выпуклое множество. Требуем еще, чтобы оно было непустым, тогда задача (УШ. I. 8) будет иметь одно и только одно решение, которое будем рассматривать как приближенное решение задачи (УШ. I. 6). Известно, что если $x_0 \notin M$, то при n достаточно большом будет $x_{on} \in M$ и, по доказанному в п. 1°, $x_*^{(n)} \in \partial M$.

5°. Известны многие способы решения задач типа (Уш. I.8); здесь ограничимся ссылкой на гл. II книги Гловинского, Дюкса и Тремольера [1]. Укажем еще один прием, который может оказаться бесполезным: см. статью автора [27].

Пусть $y_{n1}, y_{n2}, \dots, y_{nN}$ — базис в H_n . Допустим, что существует такая монотонная функция $\omega: R_+^N \rightarrow R_+^N$, что

$$\omega(\|x^{(n)} - g\|) \quad , \quad \text{где} \quad x^{(n)} = \sum_{k=1}^N a_k y_{nk}$$

здесь достаточно гладкая функция коэффициентов a_k . Положим

$$x_k^{(n)} = \sum_{k=1}^N a_k^{(n)} y_{nk} ;$$

обозначим через $a^{(n)}$ вектор $(a_1^{(n)}, a_2^{(n)}, \dots, a_N^{(n)})$, через $M_n^{(n)}$ — матрицу чисел $[y_{nk}, y_{nj}]$; $j, k = 1, 2, \dots, N$. Положим еще $C_j^{(n)} = [u_{0n}, y_n]$, $C^{(n)} = (C_1^{(n)}, C_2^{(n)}, \dots, C_N^{(n)})$ и $\Phi_n(a^{(n)}) = \omega(\|x_k^{(n)} - g\|)$. По способу множителей Лагранжа задача (Уш. I.8) сводится к решению уравнений относительно коэффициентов $a_j^{(n)}$:

$$\frac{\partial}{\partial a_j^{(n)}} (M_n a^{(n)}, a^{(n)}) - 2\lambda \frac{\partial \Phi_n(a^{(n)})}{\partial a_j^{(n)}} = 2C_j^{(n)}, \quad 1 \leq j \leq N,$$

или, короче,

$$M_n a^{(n)} - \lambda g \text{ grad } \Phi_n(a^{(n)}) = C^{(n)}. \quad (\text{Уш. I.9})$$

Будем рассматривать λ как независимую переменную и $a_j^{(n)}$ — функции от λ . Дифференцируя (Уш. I.9) по λ , получим систему обыкновенных дифференциальных уравнений

$$(M_n - \lambda \Psi_n) \frac{da^{(n)}}{d\lambda} - g \text{ grad } \Phi_n(a^{(n)}) = C, \quad (\text{Уш. I.10})$$

к которой надо добавить при начальном условии

$$C^{(n)}|_{\lambda=0} = M_n^{-1} C^{(n)}; \quad (\text{Уш. I.11})$$

в уравнении Ш. I.10) Ψ_n означает матрицу вторых производных функции Φ_n .

Задачу (Уш. I.10–11) будем решать каким-нибудь разностным методом. Пусть h — шаг соответствующей сетки. Положим $\lambda_k = \pm kh$, $k = 1, 2, \dots$ пределим вектор $a^{(n)}(k)$, соответствующий

некоторому λ_k , вычислим величину $\| \sum_{j=1}^N a_j^{(n)}(k) y_{nj} - q \|$. При малых K она близка к $\| u_0 - q \|$ и, следовательно, больше, чем d . Процесс решения задачи (УШ.1.10-II) следует приостановить при том значении K , положительном или отрицательном, при котором в пределах принятой точности будет выполняться равенство

$$\| \sum_{j=1}^N a_j^{(n)}(k) y_{nj} - q \| = d. \quad (\text{УШ.1.12})$$

Вычисление величины $a_j^{(n)}$ упрощается, если полунорма $\| \cdot \|$ определяется равенством $\| u \|^2 = b(u, u)$, где $b(u, v)$ - симметричная билинейная форма, определенная на всем пространстве H . Если в этом случае положить $\omega^{(+)} = t^2$, то система (УШ.1.9) примет вид

$$(M_n - \lambda F_n) a^{(n)} = c^{(n)} - f^{(n)}, \quad (\text{УШ.1.13})$$

где F_n - матрица чисел $b(y_{nj}, y_{nk})$; $j, k = 1, 2, \dots, N$, а $f^{(n)}$ - вектор (составляющие $f_k^{(n)} = b(q, y_{nk})$). Матрица M_n положительно определенная и потому невырожденная; систему (УШ.1.13) можно решить, обратив матрицу $M_n - \lambda F_n$ по методу Гаусса (Гаусс и Гаусса [1]; см. также статью Вайнштейна и автора [1]). Определив коэффициенты $a_j^{(n)}$, мы затем найдем λ из уравнения $b(x^{(n)} - q, x^{(n)} - q) = d^2$.

§ 2. Точность аппроксимации.

Г⁰. Пусть B - банахово пространство с нормой $\| \cdot \|$, обладающее следующими свойствами:

1) Γ гомоморфно вкладывается в \tilde{H} ; при этом

$$\forall x \in B, |x| \leq \|x\|, \|x\| \leq \|x\|. \quad (\text{УШ.2.1})$$

2) элементы α_0 и α_* содержатся в B ;

3) при любом n подпространство $H_n \subset B$, и последовательность $\{H_n\}$ полна в B .

Если в силу а) $\|\alpha_0\|$ конечна, что мы ниже всегда предполагаем, то условиям 1)-2) можно удовлетворять, положив, например,

$$\|x\|^q = |x|^q + \|x\|^q, 1 \leq q < \infty; \quad (\text{УШ.2.2})$$

дадим класс случаев, когда условие 3) также выполняется. Допустим, например, что $\tilde{H} \sim W_2^{(K)}(\Omega)$, где K - натуральное число и Ω - конечная область евклидова пространства, удовлетворяющая усло-

вну конуса, и пусть для некоторых $\delta \geq k \geq p \geq 1$ справедливо неравенство $\| \cdot \| \leq C \| \cdot \|_{W_p^{(\delta)}(\Omega)}$, $C = \text{const}$. В качестве H_n возьмем подпространства, натянутые на координатные функции МКЭ (см. книгу автора [22]), которые в свою очередь получены преобразованием независимых переменных из соответствующих функций; эти последние можно выбрать так, чтобы последовательность $\{H_n\}$ была полной в $W_p^{(\delta)}(\Omega)$. Пусть еще p таково, что $W_p^{(\delta)}(\Omega)$ ограниченно вкладывается в $W_2^{(k)}(\Omega)$. Тогда $H_r \subset V$ для любого n и последовательность $\{H_n\}$ полна в V .

Ниже мы будем пользоваться равенством (УШ.2.2) при $q = 1$.

2°. Лемма УШ.2.1. Пусть M - выпуклое замкнутое в норме $|\cdot|$ непустое множество. Пусть, далее, x_* - решение односторонней вариационной задачи

$$|x_0 - x_*| = \min, \quad x_0 \in N \setminus M, \quad x_* \in M, \quad (\text{УШ.2.3})$$

а x - произвольный элемент множества M . Тогда

$$|x_* - x|^2 \leq d'^2 - d^2, \quad (\text{УШ.2.4})$$

$$d = |x_0 - x_*|, \quad d' = |x_0 - x|. \quad (\text{УШ.2.5})$$

В пространстве H проведем двумерную плоскость P через точки x_0, x_*, x . Гильбертова метрика пространства H индуцирует на P евклидову метрику. Построим треугольник с вершинами x_0, x_*, x и обозначим через ψ угол при вершине x_* . Докажем, что $\psi \geq \pi/2$. Допустим, что $\psi < \pi/2$. Через точку x_0 проведем в P прямую, наклонную к прямой $\overline{x_0 x_*}$ под углом, меньшим, чем $\pi/2 - \psi$, и пересекающую отрезок $\overline{x_* x}$; Пусть y - точка пересечения (рис.1). В треугольнике с вершинами x_0, x_*, y угол при вершине y тупой, поэтому $|x_0 - y| < d$. Множество M выпуклое, а $x, x_* \in M$, поэтому $y \in M$. Теперь неравенство $|x_0 - y| < d$ означает, что x_* не есть решение задачи (УШ.2.3), вопреки предположению. Таким образом, $\psi \geq \pi/2$.

Теперь по теореме косинусов,

$$d'^2 = |x_* - x|^2 + d^2 - 2d|x_* - x|\cos\psi,$$

по $\cos \psi \leq 0$, и $d'^2 \geq |x_* - x|^2 + d^2$. Неравенство (УШ.2.4) доказано.

3^o. Лемма УШ.2.2. Пусть выполнены условия 1) - 3) п. 1^o настоящего параграфа, множество ограничений M определено формулой (УШ.1.1), причем $\beta = \|g\| < d$. Пусть x - некоторая точка на ∂M , так что $\|x - g\| = d$. Допустим, что существует точка $x^{(n)} \in H_n$, удовлетворяющая неравенству $\|x - x^{(n)}\| \leq r_n$, где $r_n < (d - \beta)/2$, и $\|\cdot\| = |\cdot| + \|\cdot\|$. Тогда существует элемент $y^{(n)} \in M \cap H_n$, такой, что $\|x - y^{(n)}\| \leq C_x r_n$, где C_x может зависеть от x , но не зависит от n .

Положим $y^{(n)} = ax^{(n)}$, $0 < a < 1$; очевидно, $y^{(n)} \in H_n$. Выясним, при каких a справедливо включение $y^{(n)} \in M$. Имеем

$$\|y^{(n)} - g\| \leq a \|x^{(n)} - g\| + (1-a)\beta,$$

$$\|x^{(n)} - g\| \leq \|x^{(n)} - x\| + \|x - g\| \leq r_n + d.$$

Отсюда $\|y^{(n)} - g\| \leq a(r_n + d) + (1-a)\beta$. Если

$a(r_n + d) + (1-a)\beta = d$, то

$$a = \frac{d - \beta}{d - \beta + r_n} \quad (\text{УШ.2.6})$$

то $y^{(n)} \in M$. При этом

$$\|x - y^{(n)}\| \leq \|x - x^{(n)}\| + (1-a)\|x^{(n)}\| \leq r_n + \frac{r_n \|x^{(n)}\|}{d - \beta + r_n}.$$

Далее, $\|x^{(n)}\| \leq \|x\| + \|x - x^{(n)}\| \leq \|x\| + r_n$, следовательно,

$$\|x - y^{(n)}\| \leq r_n \left[1 + \frac{\|x\| + r_n}{d - \beta + r_n} \right]. \quad (\text{УШ.2.7})$$

Отсюда вытекает более простая оценка

$$\|x - y^{(n)}\| \leq C_x r_n, \quad C_x = 2 \left[1 + \frac{\|x\|}{d - \beta} \right]. \quad (\text{УШ.2.8})$$

Лемма доказана. Аналогично можно получить оценки

$$\|x - y^{(n)}\| \leq C'_x r_n, \quad \|x - y^{(n)}\| \leq C''_x r_n, \quad (\text{Л.2.9})$$

где

$$C'_x = 2 \left[1 + \frac{\|x\|}{d - \beta} \right], \quad C''_x = 2 \left[1 + \frac{\|x\|}{d - \beta} \right]. \quad (\text{Л.2.10})$$

Теорема УШ.2.1. Если пространство B удовлетворяет условиям 1) - 3) п. I⁰, а $x_0, x_1, x_n^{(n)}$ имеют тот же смысл, что и выше, то

$$\|x_n - x^{(n)}\| = O(\sqrt{e_n(x_1)}), \quad (\text{УШ.2.11})$$

где $e_n(x)$ есть наилучшее приближение элемента $x \in B$ элементами подпространства H_n в норме $\|\cdot\|$.

Заметим, что по предположению последовательность $\{H_n\}$ полна в B , поэтому $e_n(x) \xrightarrow{n \rightarrow \infty} 0$. В частности, $e_n(x_1) \xrightarrow{n \rightarrow \infty} 0$. Как было доказано в п.4⁰ § I, $x_n^{(n)} \in \mathcal{M}$, если n достаточно велико.

По определению наилучшего приближения, при любом заданном $\varepsilon > 0$ существует такой элемент $x^{(n)} \in H_n$, что $\|x_n - x^{(n)}\| \leq (1 + \varepsilon)e_n(x_1)$. По лемме УШ.2. существует элемент $y^{(n)} \in M_n = M$ так что

$$\|x_n - y^{(n)}\| \leq C_* e_n(x_1), \quad C_* = (1 + \varepsilon)C_{x_1}.$$

$$\text{По лемме УШ.2 I, } \|x_n - x_n^{(n)}\| \leq d_n^2 - d^2, \quad d_n^2 = \|x_0 - x_n^{(n)}\|^2.$$

Но $x_n^{(n)}$ есть решение задачи УШ.2.7), поэтому $d_n \leq d'_n = \|x_0 - y^{(n)}\|$.

$$\|x_n - x_n^{(n)}\| \leq d_n^2 - d^2 = (d'_n + d)(d'_n - d).$$

Кроме того,

$$d'_n = \|x_0 - y^{(n)}\| \leq \|x_0 - x_n\| + \|x_n - y^{(n)}\| \leq d + C_* e_n(x_1).$$

Отсюда

$$\|x_n - x_n^{(n)}\|^2 \leq [2d + C_* e_n(x_1)] C_* e_n(x_1). \quad (\text{УШ.2.12})$$

Теорема доказана

§ 3. Случай улучшения оценки погрешности аппроксимации.

Существуют случаи, когда оценка (УШ.2.11) может быть существенно улучшена. В настоящем параграфе приведены два таких случая описанных в статьях автора [27] и [28].

В п. I заменим точку x на $x_n^{(n)}$. По теореме косинусов

$$\|x_n - x_n^{(n)}\|^2 = \|x_0 - x_n^{(n)} - d\|^2 + 4d \|x_0 - x_n^{(n)}\| \sin^2 \frac{\theta}{2};$$

через θ здесь обозначен угол при вершине x_0 в треугольнике $(x_0, x_n, x_n^{(n)})$. Поскольку $\|x_0 - x_n^{(n)}\| \xrightarrow{n \rightarrow \infty} d$, то при n достаточно большом $\|x_0 - x_n^{(n)}\| \leq 2d$. Обозначим еще через $x^{(n)}$ элемент подпространства H_n , на котором достигается равенство

$|x_* - x^{(n)}| = e_n(x_*)$, тогда

$$d \leq |x_0 - x_*^{(n)}| \leq |x_0 - x^{(n)}| \leq |x_0 - x| + |x - x^{(n)}| = d + e_n(x_*);$$

отсюда $0 \leq |x_0 - x_*^{(n)}| - d \leq e_n(x_*)$ и, следовательно,

$$|x_* - x_*^{(n)}|^2 \leq e_n^2(x_*) + 2d^2 \theta_n^2, \quad (\text{Лш. 1})$$

и дело сводится к оценке угла θ_n . Допустим, что многообразие \mathcal{M} достаточно гладко, тогда вектор $x_0 - x_*$ направлен по нормали к \mathcal{M} . Приведем касательство этого достоящего утверждения. Выберем на \mathcal{M} произвольную точку y и проведем через точки x_0, x_*, y поперечную плоскость Q (Лш. 2), которая пересечет \mathcal{M} по достаточно гладкой кривой \mathcal{L} . Введем в Q декартову систему координат (ξ_1, ξ_2) ; за ось ξ_1 примем прямую, проходящую через точки x_0 и x_* и направленную от x_0 к x_* ; за ось ξ_2 примем за ось полярных координат, центр которых поместим в x_0 . Уравнение касательной к кривой \mathcal{L} в точке x_* можно написать в виде $(\xi_1 - \xi_{*1} - \text{осциллыса переменной точки } x)$

$$\frac{\xi_1 - \xi_{*1}}{d\xi_1} = \frac{\xi_2}{d\xi_2}.$$

Пусть $\rho = \rho(\theta)$ есть уравнение кривой \mathcal{L} в полярных координатах вблизи точки x_* , тогда

$$d\xi_{*1} = d(\rho \cos \theta) = \cos \theta d\rho - \rho \sin \theta d\theta.$$

Но в точке x_* $d\rho = 0$, потому что в этой точке ρ достигает минимума; кроме того, в этой же точке $\sin \theta = 0$. Таким образом $d\xi_{*1} = 0$ и угловая касательная к \mathcal{L} в точке x_* есть $\xi_2 = \xi_{*2}$, и прямая x_*x_0 направлена по нормали к \mathcal{M} . Как было только-что отмечено, в введенных нами полярных координатах можно записать уравнение кривой \mathcal{L} вблизи точки x_* в виде $\rho = \rho(\theta)$, где $\rho(\theta)$ - достаточно гладкая функция от θ , причем $\rho(0) = d$, $\rho'(0) = 0$. Наложим на многообразие \mathcal{M} следующее требование T:

Существуют такие положительные числа K_0, K_1, \bar{K} , что

$$K_0 \leq \frac{1}{2} \rho''(0) \leq K_1; \quad \frac{1}{6} |\rho'''(\theta)| \leq \bar{K}; \quad (\text{Лш. 3.2})$$

эти числа могут зависеть от X_0 , но не зависят от выбора двумерного сечения, проходящего через точки X_0 и X_* ; второе равенство (УШ.3.2) должно выполняться при достаточно малых θ .

Обозначим $\rho(\theta) - d = l(\theta)$. Разлагая $\rho(\theta)$ по формуле Тейлора, находим $l(\theta) = k\theta^2 + \rho'''(\tau)\theta^3/6$; $k = \rho''(0)/2$.

Отсюда $\theta^2 = l(\theta)/k - \rho'''(\tau)\theta^3/6k$ и, следовательно,

$$\frac{1}{k_1} l(\theta) - \frac{\bar{K}}{k_0} |\theta|^3 \leq \theta^2 \leq \frac{1}{k_0} l(\theta) + \frac{\bar{K}}{k_0} |\theta|^3.$$

Пусть $|\theta| \leq k_0/2\bar{K}$. Из последнего неравенства находим

$$\frac{2}{3k_1} l(\theta) \leq \theta^2 \leq \frac{2}{k_0} l(\theta). \quad (\text{УШ.3.3})$$

По теореме косинусов

$$\begin{aligned} |x_* - y|^2 &= d^2 + \rho^2(\theta) - 2d\rho(\theta)\cos\theta = \\ &= \rho^2(\theta) + 4d(d + l(\theta))\sin^2\frac{\theta}{2}; \end{aligned}$$

отсюда и из неравенства (УШ.3.3) вытекает неравенство

$$\beta_1 \theta \leq |y - x_*| \leq \beta_2 \theta \quad (\text{УШ.3.4})$$

с достоящими β_1 и β_2 , которые не зависят от выбора описанного выше двумерного сечения.

Теорема УШ.3.1. Пусть выполнено требование Т, а также предположены п.1^o § 2. Тогда

$$|x_* - x_*^{(n)}| \leq C \epsilon_n(x_*), \quad C - \text{const}. \quad (\text{УШ.3.5})$$

Существует такая n элемент $x_*^{(n)} \in U_n$, что $|x_* - x_*^{(n)}| \leq (1 + \epsilon) \epsilon_n(x_*)$. По лемме УШ.2.2 найдется элемент $y^{(n)} \in M_n$ удовлетворяющий неравенству $|x_* - y^{(n)}| \leq C_* \epsilon_n(x_*)$. Воспользуемся обозначениями рис.3; все расстояния и углы измеряются в радианах.

По построению приближенного решения, $\rho'_n \leq \rho'_n$. Обозначим $\epsilon'_n = \rho'_n - \rho'_n$ и $l'_n = \rho'_n - d$, тогда $l'_n < l'_n$.

в силу неравенства (Уш.3.3)

$$\theta_n^2 \leq \frac{2}{K_1} l_n < \frac{2}{K_0} l'_n \leq \frac{2K_1}{K_0} \theta_n'^2,$$

или $|\theta_n| < |\theta'_n| \sqrt{3K_1/K_0}$. Теперь по неравенству (Уш.3.4),

$$|x_* - x'_*| \leq \beta_2 \theta_n < \beta_2 \theta'_n \sqrt{3K_1/K_0} \leq$$

$$\leq \frac{\beta_2}{\beta_1} \sqrt{3K_1/K_0} |x_* - y^{(n)}| \leq \frac{\beta}{\beta_1} \sqrt{3K_1/K_0} C_* e_n(x_*).$$

Таким образом, если выполнено требование Т, то погрешность аппроксимации приближенного решения есть величина порядка $O(e_n(x_*))$.

2°. В настоящем пункте мы укажем некоторый простой класс функций, удовлетворяющих требованию Т. Именно, пусть $\| \| X \| \|^2 = b(x) = b(x, x)$, где $b(x, y)$ — симметричная билинейная форма, заданная на H , к которой соответствует неотрицательная квадратичная форма $b(x)$. Для упрощения выкладок будем в данном пункте считать, что $q=0$, $H = \mathbb{R}^n$.

Произвольную двумерную плоскость, проходящую через точки x_0 и x_* , можно построить так. В пространстве H возьмем такую точку y , что $|x_0 - y| = d$ и $[x_0 - x_*, x_0 - y] = 0$, и положим $e_1 = (x_* - x_0)/d$, $e_2 = (y - x_0)/d$. Двумерное множество, определенное уравнением $x - x_0 = \xi e_1 + \eta e_2$, где ξ и η — произвольные вещественные числа, есть плоскость, проходящая через точки x_0, x_*, y ; элементы вида $x - x_0$, соответствующие числам $\xi=0$ и $\eta=0$, ортогональны, поэтому ξ и η можно рассматривать как декартовы координаты на этой плоскости; начальной координат соответствующей точке x_0 .

Введем полярные координаты по обычным формулам $\xi = \rho \cos \theta$, $\eta = \rho \sin \theta$. Отметим полярные координаты точек x_0, x_*, y ; именно, $x_0: \rho=0$; $x_*: \rho=d, \theta=0$; $y: \rho=d, \theta=\pi/2$. В полярных координатах $x = x_0 + \rho u_\theta/d$, где $u_\theta = (x_* - x_0) \cos \theta + (y - x_0) \sin \theta$.

Если $\omega(x, y)$ — симметричная билинейная форма, то $\omega(x) := \omega(x, x) = d^2 \rho^2 \omega(u_\theta) + 2d \rho \omega(x_0, u_\theta) + \omega(x_0)$.

В частности, многообразия ∂M не делится уравнением

$$d^{-2} \rho^2 b(v_0) + 2d^{-1} \rho b(x_0, v_0) + b(x_0) - d = 0. \quad (\text{УШ.3.5})$$

Нетрудно видеть, что $b(v_0) > 0$. Действительно, если

$b(v_0) = 0$, то $\|x_0 - x_*\| = 0$ и $\|x_0\| = \|x_*\| = d$; это означает, что v_0 есть решение односторонней вариационной задачи вопреки предположению.

Итак, $b(v_0) > 0$. Что тогда при достаточно малых θ и $b(v_\theta) > 0$. Уравнение (УШ.3.7) — квадратное относительно ρ , из его корней один равен $\rho(\theta)$, а другой больше, чем $\rho(\theta)$. Отсюда вытекают неравенства

$$-b(x_0, v_\theta) - b(x_0, x_0 - x_*) > 0, \quad (\text{УШ.3.6})$$

$$b(x_0) - d^2 > 0. \quad (\text{УШ.3.7})$$

Положив в (УШ.3.7) $\theta = 0$, $\rho = d$, получим

$$b(x_0) - d^2 = -[f(v_0) + 2b(x_0, v_0)] = b(x_0 - x_*, x_0 - x_*),$$

и уравнение (УШ.3.7) можно записать в виде

$$d^{-2} \rho^2 b(v_\theta) + 2d^{-1} \rho b(x_0, v_\theta) + b(x_0 + x_*, x_0 - x_*) = 0, \quad (\text{УШ.3.8})$$

и неравенство (УШ.3.9) — в виде

$$b(x_0 + x_*, x_0 - x_*) > 0. \quad (\text{УШ.3.9})$$

Вычислим теперь величину $\rho''(0)$. Дифференцируем тождество (УШ.3.7) два раза по θ . Положив $\theta = 0$, получим

$$\rho''(0) = d \frac{b(x_0, x_0 - x_*) + b(x_0 - x_*)}{b(x_0, x_0 - x_*)}. \quad (\text{УШ.3.10})$$

Далее, что знаменатель дроби (УШ.3.10) положителен. При $\theta = 0$ уравнение (УШ.3.10) принимает вид

$$d^{-2} \rho^2 b(x_0 - x_*) - 2d^{-1} \rho b(x_0, x_0 - x_*) + b(x_0 + x_*, x_0 - x_*) = 0; \quad (\text{УШ.3.11})$$

его корни суть

$$\frac{d}{b(x_0 - x_*)} \left\{ b(x_0, x_0 - x_*) \left[b^2(x_0, x_0 - x_*) - b(x_0 - x_*) b(x_0 + x_*, x_0 - x_*) \right]^{1/2} \right\} =$$

$$= \frac{d [b(x_0, x_0 - x_*) - b(x_*, x_0 - x_*)]}{b(x_0 - x_*)} \quad (\text{УШ.З.14})$$

Если в (УШ.З.14) выбрать знак минус, то получим меньший корень уравнения (УШ.З.13) ρ' ; отсюда следует, что $b(x_*, x_0 - x_*) > 0$.

Из формулы (УШ.З.12) следует теперь, что $\rho''(\theta) \geq d$. С другой стороны, так как форма $O(x)$ задана на всем пространстве H , существует такая постоянная $\bar{c} > 0$, что $b(x_0 - y) \leq \bar{c} |x_0 - y|^2 = \bar{c} d^2$, и мы приходим к оценке

$$K_0 := \frac{d}{2} \leq \frac{1}{2} \rho''(0) \leq \frac{d}{2} + \frac{\bar{c} d^3}{2b(x_*, x_0 - x_*)} := K_1, \quad (\text{УШ.З.15})$$

в которой K_0 и K_1 не зависят от выбора элемента y .

Докажем теперь, что при малых θ третья производная $\rho'''(\theta)$ ограничена по модулю независимо от ρ и y . Меньший корень уравнения (УШ.З.11) равен

$$\mu(\theta) = \frac{d}{b(v_\theta)} [b(x_0, -v_\theta) -$$

$$-\sqrt{b^2(x_0, -v_\theta) - b(x_0 + x_*, x_0 - x_*) b(v_\theta)}]. \quad (\text{УШ.З.16})$$

Как мы видели, $b(v_\theta) > 0$ при малых θ . Из формулы (УШ.З.16) видно, что любая производная от $\rho(\theta)$ будет ограничена, если ограничено окончательное выражение. Было бы установлено, что при $\theta = 0$ оно равно $b^*(x_*, x_0, x_*) > 0$. Далее,

$$v_\theta = x_0 - x_* - 2(x_0 - x_*) \sin^2 \frac{\theta}{2} + (x_0 - y) \sin \theta := x_0 - x_* + \eta;$$

заметим, что $|\eta| \leq d(|\theta| + \frac{\theta^2}{2})$ и потому при малых θ будет

$|\eta| \leq \beta |\theta|$, где β не зависит от θ от y . Подкоренное

выражения в (УШ.3.16) равно

$$b^2(x_0, x_0 - x_*) + b^2(x_0, \eta) b(x_0, x_0 - x_*) + b^2(x_0, \eta) - \\ - [2b(x_0 - x_*, \eta) + b(\eta)] b(x_0 + x_*, x_0 - x_*) = \\ = b^2(x_0, x_0 - x_*) + O(|\theta|),$$

и при достаточо малых θ оно ограничено снизу положительным числом, не зависящим от θ и ψ .

Заметим еще, что требование Т выполнено, если для любого элемента $y \in \mathcal{L}$ для \mathcal{L} - прямая (или ее часть) - в этом случае

$$\rho(\vartheta) = d/\cos \theta \text{ и можно поло. т.ч. } K_0 = K_1 = d/2 \text{ и, при } |\theta| \leq \pi/4 \\ K = 4d/\theta.$$

°. Как и выше, возьмем произвольную точку $y \in \mathcal{M}$ и проведем в \mathbb{H} двумерную плоскость P через точки x_0, x_*, y (см. рис. 2) и введем на этой плоскости полярные координаты ρ и θ . Отметим, что полунорма $\| \cdot \|$ конечна на всех элементах плоскости P ; кроме того, как было показано в § 2, $\psi \geq \pi/2$. Пусть по-прежнему $L = \mathcal{M} \cap P$ и $\rho(\theta) = \rho(\theta)$ есть уравнение кривой L . Функция $\rho(\theta)$ достигает минимума при $\theta = 0$ поэтому $\rho(\theta)$ возрастает при возрастании $|\theta|$. Из выпуклости кривой L следует, что с возрастанием θ угол ψ также возрастает.

Сделаем теперь следующее допущение: существуют такие положительные постоянные $\theta_0, k, \delta < \pi/2$, что при $|\theta| \leq \theta_0$ справедливы неравенства $|\rho'(\theta)| \leq k$ и $\delta + \pi/2 \leq \psi$. Эти постоянные могут зависеть от x_0 , но не от выбора двумерной плоскости P , проходящей через точки x_0 и x_* .

Из неравенства $\psi \geq \delta + \pi/2$ следует, что x_* - "конечная" точка множества \mathcal{M} : вектор $x_0 - x_*$ образует угол, не меньший чем $\delta + \pi/2$, с любой секущей множества \mathcal{M} , проходящей через точку x_* .

Покажем, что при сформулированном здесь допущении верна оценка (УШ.3.5) для погрешности приближенного решения.

Пусть $|\theta| \leq \theta_0$. Обозначим $\ell(\theta) = \rho(\theta) - d$. Имеем

$$\rho(\theta) = \rho(\theta) \cdot 1 = \rho(\theta) \cdot \rho(0) = \int_0^\theta \rho(\vartheta) d\vartheta \leq k|\theta|, \quad (\text{УШ.3.4})$$

поэтому, если x - точка с полярными координатами $(\rho(\theta), \theta)$, то

$$|\alpha_* - \alpha|^2 = \rho^2(\theta) + d^2 - 2\rho(\theta)d \cdot \cos \theta =$$

$$= \ell^2(\theta) + 4\rho(\theta)d \cdot \sin^2 \frac{\theta}{2} \leq C_1^2 \theta^2, \quad C_1 = \text{const.} \quad (\text{УШ.3.18})$$

С другой стороны, обозначая $\lambda = |\alpha_* - \alpha|$, находим

$$\lambda \leq |\theta| \frac{\lambda \sin \psi}{\rho(\theta)}; \quad |\theta| \leq \frac{\theta_0}{\sin \theta_0} \frac{\lambda}{d} =: C_2 |\alpha - \alpha_*|. \quad (\text{УШ.3.19})$$

Кроме того,

$$\lambda = \frac{\rho(\theta) \sin \theta}{\sin \psi} \geq \frac{\rho(\theta) \sin \theta}{\cos \delta}$$

и, следовательно,

$$\ell^2(\theta) + 4\rho(\theta)d \cdot \sin^2 \frac{\theta}{2} \geq \frac{\rho^2(\theta) \sin^2 \theta}{\cos^2 \delta}. \quad (\text{УШ.3.20})$$

Докажем, что при достаточно малом θ_0 справедливо неравенство

$$4\rho(\theta)d \cdot \sin^2 \frac{\theta}{2} \leq \sigma \rho^2(\theta) \frac{\sin^2 \theta}{\cos^2 \delta}, \quad (\text{УШ.3.21})$$

в котором σ — г-установка, $0 < \sigma < 1$. Неравенство (УШ.3.21) равносильно следующему: $d \leq \sigma \rho(\theta) \cos^2 \frac{\theta}{2} / \cos^2 \delta$, и достаточно, чтобы $d \leq \sigma \rho(\theta) \cos^2 \frac{\theta_0}{2} / \cos^2 \delta$. Если взять $C_2 < 2\delta$, то последнее неравенство будет выполнено при $\sigma = \cos^2 \delta / \cos^2 \frac{\theta_0}{2}$.

Теперь

$$\ell^2(\theta) \geq (1-\sigma) \frac{\rho^2(\theta) \sin^2 \theta}{\cos^2 \delta} > (1-\sigma)d \frac{\sin^2 \theta}{\cos^2 \delta},$$

и окончательно

$$\ell(\theta) \geq C_3 |\theta|, \quad C_3 = \frac{d \sqrt{1-\sigma} \sin \theta_0}{\theta_0 \cos \delta}. \quad (\text{УШ.3.22})$$

Пусть $\bar{x}^{(n)}$ элемент подпространства H , на котором достигнута наилучшая приближение элемента x_* , т.е. что $E_n(x_*) = \|x_* - \bar{x}^{(n)}\|$. По лемме УШ.2. существует элемент $y^{(n)} \in M \cap H_n$, такой, что $\|x_* - y^{(n)}\| \leq C E_n(x_*)$. Воспользуемся обозначениями рис.3, в котором на этот раз будем считать $x^{(n)} = y^{(n)}$; все углы и расстояния измеряются в метрике l . Построим

приближенного решения, $\rho'_n > \rho_n$. Обозначим $l_n = \rho_n - d$, $l'_n = \rho'_n - d$, так что $l_n \cdot l'_n$. По формулам (УШ.З.18), (УШ.З.22) (УШ.З.17) - (УШ.З.19) получаем

$$|x_n - x_n^{(m)}| \leq C_1 |\theta_n| \leq \frac{C_1}{C_3} l_n < \frac{C_1}{C_3} l'_n < \frac{C_1 k}{C_3} |\theta'_n| \leq \\ \leq \frac{C_1 C_2 k}{C_3} |x_n - y^{(m)}| \leq C e_n(x_n); C = \frac{C_1 C_2 k}{C_3},$$

и наше утверждение доказано.

§ 3. Оценка точности аппроксимации для более общей задачи

Г. В работе Тохтина [1] исследована точность аппроксимации более общей задачи, сформулированной в начале настоящей главы. Напомним постановку этой задачи. В рефлексивном банаховом пространстве B с нормой $\|\cdot\|$ заданы непустое выпуклое замкнутое множество M и выуклый слабо непрерывный оператор J ; если множество M неограниченное, то дополнительно потребуем, чтобы $J(x) \xrightarrow{\|x\| \rightarrow +\infty} +\infty$. Ставится задача о нахождении элемента $x_n \in M$, реализующего минимум функционала $J(x)$ на множестве M . Как было сказано в начале главы, решение этой задачи существует и единственно, если функционал J строго выуклый.

2°. Приближенное решение $x_n^{(n)}$ этой более общей задачи с тем же построением так же, как в § 2: выберем полную в B последовательность конечномерных подпространств $\{B_n\}$ и примем за $x_n^{(n)}$ решение задачи о минимуме функционала J на множестве $M \cap B_n$. Поскольку последовательность $\{B_n\}$ в цитируемой статье Тохтина построена в следующем смысле. Пусть $R_M(x)$ функционал минковского для множества M . Этот функционал определяется следующим образом. Допустим, что L - произвольное линейное пространство и \mathcal{M} - выуклое множество в L . Точка $\xi \in M$ принадлежит внутренней множеству M (термин взят из книги Демидова и Вомина [1]; см. также Баланшина [1]) называется алгебраически внутренней точкой множества M , если для любого $\eta \in L$ существует такое $\alpha > 0$ что $(\xi + \alpha \eta) \in M$, а $|\alpha|^{1+\epsilon} < \alpha$. Пусть ξ_0 - точка ядра множества M . Функция на минковского множества M определяется формулой (см. цитированные выше книги

$$\forall x \in L, P_M(x) = \inf_{\lambda > 0} \{ \lambda : \frac{x - \xi}{\lambda} \in M \}. \quad (Уш.4.1)$$

Ниже будем считать, что $\xi_0 = 0$ - это, очевидно, не ограничивает общности.

Положим теперь

$$\forall x \in B, \varepsilon_n(x) = \inf_{x' \in B_n} [\|x - x'\| + |P_M(x) - P_M(x_n)|]. \quad (Уш.4.2)$$

Будем считать, что на элементах подпространств B_n функционал Линковского конечен. По определению, поэтому в работе Тихтина [1], последовательность $\{B_n\}$ полна в точке $x \in B$, если

$$\lim_{n \rightarrow \infty} \varepsilon_n(x) = 0. \quad \text{Очевидно, для такой полноты необходимо, чтобы } P_M(x) < \infty.$$

Введем в рассмотрение функцию

$$q(\alpha) = \inf_{y \in M_\alpha} [\gamma(y) - \gamma(x_*)], \quad (Уш.4.3)$$

$$M_\alpha = \{ y \in M; \|y - x_*\| = \alpha \},$$

которую назовем инерантой функционала γ на множестве M .

Оценка погрешности аппроксимации дается следующей теоремой:

Теорема Уш.4.1. Пусть выполнены следующие условия: а) на всех элементах всех подпространств B_n функционал Линковского конечен; б) в некоторой окрестности x_* (в смысле нормы пространства R) точного решения x_* функционал $\gamma(x)$ конечен и удовлетворяет условию Липшица; в) при малых значениях числа $\alpha \geq 0$ функция $q(\alpha)$ строго монотонна; г) последовательность подпространств $\{B_n\}$ полна в B в указанном выше смысле. Тогда имеет место оценка погрешности аппроксимации

$$\|x_n - x_*^{(n)}\| \leq q^{-1}(\beta \varepsilon_n(x_*)); \quad \beta = \beta \max(1, \|x_*\|), \quad (1.4)$$

β - постоянная Липшица для функционала γ в окрестности точного решения x_* .

Замечание Уш.4.1. В работе Тихтина [1] доказано, что γ удовлетворяет условию Липшица, если M имеет внутреннюю точку.

Замечание Уш.4.1. При доказательстве теоремы Уш.4.1 играет

важную роль лемма, которая является обобщением леммы 4.2.:

Лемма УШ.4.1. Пусть $X \in M$ и $x^{(n)} \in B_n$. Существует точка

$y^{(n)} \in B_n \cap M$, такая, что

$$\|x - y^{(n)}\| \leq \|x - x^{(n)}\| + \|x\| \cdot |P_M(x) - P_M(x^{(n)})|. \quad (\text{УШ.4.5})$$

Отметим важный частный случай, рассмотренный в цитированной работе Лактина. Пусть B есть гильбертово пространство и $f(x) = \|x - x_0\|^2$, $x_0 \notin M$; причём, что M произвольное непустое выпуклое замкнутое множество в B . Тогда справедлива

Теорема УШ.4.2. Если функционал Минковского P_M конечен на всех элементах пространства B_n и последовательность $\{B_n\}$ полна в точке x_* в упомянутом смысле, то

$$\|x_* - x_*^{(n)}\| = O(\sqrt{\epsilon_n(x_*)}). \quad (\text{УШ.4.6})$$

4°. В настоящем пункте мы изложим, также без доказательства, некоторые из результатов работ Демьяновича [4]. В ней рассмотрена задача п.3 § I с минимумом расстояния от заданной точки $x_0 \in H$, где H - гильбертово пространство, до непустого выпуклого замкнутого в H множества M . Пусть H_n - подпространство пространства H ; x_* , $x_*^{(n)}$ и $x^{(n)}$ имеют тот же смысл, что и во всем предыдущем. Обозначим

$$\forall x \in H, E_n(x) = \inf_{x^{(n)} \in H_n} \|x - x^{(n)}\|; \quad (\text{УШ.4.7})$$

кроме того, пусть P_n - ортогональный проектор из H на H_n . Доказывается следующая теорема:

Теорема УШ.4.3. Пусть множество $M_n = M \cap H_n$ и $P_n x_* \in M_n$. Тогда имеет место оценка

$$\|x_* - x_*^{(n)}\| \leq E_n(x_*) + [E_n(x_*) + E_n(x_*) E_n(x_0)]^{1/2}. \quad (\text{УШ.4.8})$$

Заметим, что если $E_n(x_*)$ и $E_n(x_0)$ суть величины одного порядка, то из (УШ.4.8) следует, что

$$\|x_* - x_*^{(n)}\| = O(E_n(x_*)) \quad (\text{УШ.4.9})$$

к сожалению, условие $P_n x_* \in M_n$ трудно проверить.

§ 5. Об устойчивости точек равновесия односторонней вариационной задачи

1°. Здесь рассматривать задачу о минимуме расстояния в

гильбертовом пространстве H от фиксированной точки x_0 этого пространства до непустого вышуклого замкнутого множества M :

$$\|x - x_0\| = \min, x \in M. \quad (\text{УШ.5.1})$$

поставим вопрос (см. работу автора [29]) об устойчивости точного решения задачи (УШ.5.1) относительно малых изменений ее данных: элемента x_0 , нормы $\|\cdot\|$ пространства H и множества M . Этот вопрос, который кажется нам интересным и сам по себе, поможет нам потом (см. § 6) разобраться в полноте искажения приближенного решения задачи (УШ.5.1). Эту последнюю мы будем ниже называть простейшей односторонней вариационной задачей.

2°. Пусть норма $\|\cdot\|$ и множество M остаются невозмущенными, а элемент x_0 заменен элементом $y_0 \in H$, причем $\|x_0 - y_0\| \leq \delta$; через δ здесь и ниже в данном параграфе обозначается малое положительное число. Решение возмущенной таким образом задачи (УШ.5.1) обозначим через y_* . Так как, по обычному предположению, $x_0 \notin M$, то при достаточно малом δ будет также $y_0 \notin M$, а тогда $x_* \neq x_0$, $y_* \in y_0$, $x_* \in \partial M$, $y_* \in \partial M$.

Рассмотрим четырехугольник рис.49. Его четыре вершины определяют трехмерное (в частном случае - двумерное) пространство E , все элементы которого принадлежат H . Из доказательств леммы УШ.2.1 вытекает, что углы четырехугольника в вершинах x_* и y_* - тупые или прямые, поэтому, если через точки x_* и y_* провести в E плоскости (прямые, если пространство E двумерное), нормальные к отрезку, соединяющему точки x_* и y_* , то x_0 и y_0 будут лежать вне или на границе слоя, образованного проведенными плоскостями. Отсюда следует, что

$$\|x_* - y_*\| \leq \|x_0 - y_0\| \leq \delta, \quad (\text{УШ.5.2})$$

решение простейшей задачи устойчиво относительно малых изменений элемента x_0 .

3°. В пространстве H заменим норму $\|\cdot\|$ другой гильбертовой нормой $\|\cdot\|_1$, так, чтобы

$$\forall x \in H, (1-\delta)\|x\| \leq \|x\|_1 \leq (1+\delta)\|x\|. \quad (\text{УШ.5.3})$$

пространство H , снабженное метрикой $\|\cdot\|_1$, обозначим через H_1 ,

а решение возмущенной задачи $|\mathcal{X} - \mathcal{X}_0|_1 = \min, \mathcal{X} \in \mathcal{M}$ - через \mathcal{X}_* . Оценим величину $|\mathcal{X}_* - \mathcal{X}_0|$. Через точки $\mathcal{X}_0, \mathcal{X}_*$ проведем в \mathbb{H} двумерную плоскость P и введем на ней евклидову метрику, индуцированную гильбертовой нормой $|\cdot|$. В обозначениях рис.5 имеем

$$|\mathcal{X}_0 - \mathcal{X}_*|^2 = d^2 + |\mathcal{X}_* - \mathcal{X}_1|^2 - 2d |\mathcal{X}_* - \mathcal{X}_1| \cos \Psi.$$

Очевидно, $\Psi \geq \pi/2$ - в противном случае \mathcal{X}_* не была бы точкой минимума, поэтому

$$|\mathcal{X}_* - \mathcal{X}_1|^2 \leq |\mathcal{X}_0 - \mathcal{X}_*|^2 - d^2. \quad (\text{VII.5.4})$$

Аналогично, если на плоскости P ввести евклидову метрику, индуцированную гильбертовой нормой $|\cdot|_1$, и обозначить

$|\mathcal{X}_0 - \mathcal{X}_*|_1 = d_1$, то получим

$$|\mathcal{X}_* - \mathcal{X}_1|_1^2 \leq |\mathcal{X}_0 - \mathcal{X}_*|_1^2 - d_1^2. \quad (\text{VII.5.5})$$

Если $d_1 \leq d$, то из (VII.5.5) следует

$$|\mathcal{X}_* - \mathcal{X}_1|^2 \leq \frac{|\mathcal{X}_0 - \mathcal{X}_*|_1^2}{(1-\delta)^2} - d^2 = \frac{d_1^2}{(1-\delta)^2} \cdot d^2 \leq \frac{d^2(2-\delta)}{1-\delta} \delta,$$

или

$$|\mathcal{X}_* - \mathcal{X}_1| \leq \frac{\sqrt{2-\delta} d}{1-\delta} \sqrt{\delta}; \quad (\text{VII.5.6})$$

если же $d_1 > d$, то по неравенству (VII.5.5)

$$|\mathcal{X}_* - \mathcal{X}_1|_1^2 \leq (1+\delta)d^2 - d_1^2 \leq (2+\delta)d^2\delta,$$

откуда

$$|\mathcal{X}_* - \mathcal{X}_1| \leq \frac{|\mathcal{X}_* - \mathcal{X}_1|_1}{1-\delta} \leq \frac{\sqrt{2+\delta} d}{1-\delta} \sqrt{\delta}. \quad (\text{VII.5.7})$$

Из неравенств (VII.5.6) и (VII.5.7) следует, что

$$|\mathcal{X}_* - \mathcal{X}_1| \leq C_\delta \sqrt{\delta}, \quad C_\delta = \frac{\sqrt{2+\delta}}{1-\delta} d, \quad (\text{VII.5.8})$$

и, значит, решение простейшей задачи устойчиво относительно малых возмущений метрики пространств \mathbb{H} .

4°. Допустим теперь, что элемент \mathcal{X}_0 и метрика пространств

H оставлены неизменными, а множество M заменено близким к нему множеством M' , которое мы также предполагаем непустым, выпуклым и замкнутым в норме H . Близость множеств M и M' мы здесь определяем следующим требованием: существует такое биъективное преобразование $S: M' \rightarrow M$, при котором

$$\forall x \in M, |x - S^{-1}x| \leq \delta(1 + |x|). \quad (\text{УШ.5.9})$$

Отметим одно следствие из неравенства (УШ.5.9). Пусть $x' \in M'$, тогда $Sx' \in M$ и, по неравенству (УШ.5.9)

$$|x' - Sx'| = |Sx' - S^{-1}Sx'| \leq \delta(1 + |Sx'|).$$

Полагая в том же неравенстве (УШ.5.9) $x = Sx'$, получаем $|Sx'| \leq (1 - \delta)^{-1}(|x'| + 1)$ и, следовательно,

$$\forall x' \in M', |x' - Sx'| \leq \frac{\delta}{1 - \delta}(|x'| + 1). \quad (\text{УШ.5.10})$$

Наряду с задачей (УШ.5.1) рассмотрим задачу

$$|x' - x_0| = \min, \quad x' \in M'; \quad (\text{УШ.5.11})$$

ее решение обозначим через x'_* . Пусть, как и выше, $x_0 \notin M$. Докажем, что при достаточно малом δ будет $x_0 \notin M'$. Действительно, если $x_0 \in M'$, то $Sx_0 \in M$ и $|Sx_0 - x_0| \geq d$; как и выше, здесь $d = |x_* - x_0|$. С другой стороны, по неравенству (УШ.5.10) $|x_0 - Sx_0| \leq \delta(1 + \delta)^{-1}(|x_0| + 1)$, что противоречит предшествующему неравенству, если δ достаточно мало.

Обозначим $d' = \min |x' - x_0| = |x'_* - x_0|$. По неравенству (УШ.5.9)

$$d' \leq |x_0 - S^{-1}x'_*| \leq |x_0 - x'_*| + |x'_* - S^{-1}x'_*| \leq d + \delta(1 + |x'_*|)$$

и, так как $S^{-1}x'_* \in M'$, то $d' \leq |x_0 - S^{-1}x'_*| \leq d + \delta(1 + |x'_*|)$.

Аналогично, так как $x'_* \in M'$, то $Sx'_* \in M$, и по неравенству (УШ.5.10)

$$d \leq |x_0 - Sx'_*| \leq |x_0 - x'_*| + |x'_* - Sx'_*| \leq d' + \frac{\delta}{1 - \delta}(|x'_*| + 1).$$

Таким образом,

$$-\delta(1 + |x'_*|) \leq d - d' \leq \frac{\delta}{1 - \delta}(|x'_*| + 1).$$

(УШ.5.12)

Оценим разность $x_n - x'_n$. Имеем

$$|x_n - x'_n| \leq |x_n - S^{-1}x_n| + |S^{-1}x_n - x'_n| \leq \delta(|x_n| + 1) + |S^{-1}x_n - x'_n|. \quad (\text{Ул. 5.1})$$

Точка $S^{-1}x_n \in M'$. По лемме Ул. 2.1

$$\begin{aligned} |S^{-1}x_n - x'_n| &\leq |x_0 - S^{-1}x_n| - d'^2 = \\ &= [|x_0 - S^{-1}x_n| + d] [|x_0 - S^{-1}x_n| - d'] \end{aligned}$$

и по неравенству (Ул. 5.12)

$$|S^{-1}x_n - x'_n| \leq [\delta(1-\delta)^{-1}(|x'_n| + 1) + \delta(|x_n| + 1)] \times$$

(Ул. 5.1)

$$\times [2d + \delta(1-\delta)^{-1}(|x'_n| + 1) + \delta(|x_n| + 1)].$$

Обозначим $|S^{-1}x_n - x'_n| = \varepsilon$. Из неравенства (Ул. 5.13) получаем $|x'_n| \leq (1 + \delta)(|x_n| + 1) + \varepsilon$. Подставив это в (Ул. 5.14),

приходим к квадратному неравенству для ε вида

$$(1 - a\delta^2)\varepsilon^2 - 2b\delta\varepsilon - c\delta \leq 0; \quad (\text{Ул. 5.15})$$

здесь a, b, c суть функции от δ и от $|x_n|$, ограниченные и положительные при фиксированной $|x_n|$ и при достаточно малых δ . Из (Ул. 5.15) следует, что $\varepsilon \leq C\delta^{1/2}$, $C = \text{const}$ и, в силу неравенства (Ул. 5.13),

$$|x'_n - x_n| \leq C'\sqrt{\delta}, \quad C' = \text{const}. \quad (\text{Ул. 5.16})$$

Неравенство (Ул. 5.16) показывает, что решение задачи (Ул. 5.1) устойчиво относительно возмущений множества M , малых в смысле неравенства (Ул. 5.3).

5°. Исследуем устойчивость решения задачи (Ул. 1.1-). В силу результатов п. 2^о § 1 достаточно ограничиться случаем, когда $\beta = \|g\| < d$. Устойчивость названной задачи мы исследуем в двух случаях, когда множество M' изменяется так, что это изменение естественно считать малым, но оно может быть описано соотношением вида неравенства (Ул. 5.9). В настоящем пункте мы рассмотрим случай, когда M' - возмущенное множество - опреде-

дается формулой

$$M' = \{x \in N : \|x - g\| \leq d'\}, \quad (\text{УШ.5.17})$$

в которой d' - число, близкое к d . Заметим, что если $g=0$ (и этому очевидным образом сводится более общий случай $g \in N$), то мы оказываемся в условиях п.4⁰: за S можно принять преобразование $\forall x' \in M', Sx' = dx'/d'$. Ниже рассмотрен общий случай $g \in N$.

Пусть сначала $d' = d - \delta$. Обозначим

$$M_- = \{x \in N : \|x - g\| \leq d - \delta\}. \quad (\text{УШ.5.18})$$

положим еще

$$\forall x \in M, Qx = (1 - \sigma)x = x_-, \quad \sigma > 0 \quad (\text{УШ.5.19})$$

и покажем, что при δ достаточно малом Q есть инъекция $M \rightarrow M_-$. Имеем, действительно,

$$\|x - g\| = \|(1 - \sigma)(x - g) - \sigma g\| \leq (1 - \sigma)d + \sigma\beta = d - \sigma(d - \beta)$$

достаточно принять $\sigma = \delta(d - \beta)^{-1}$; отметим еще, что $\sigma = c\delta$, $c = (d - \beta)^{-1} = \text{const}$ и что справедливо равенство

$$\|x - Qx\| = c\delta \|x\|. \quad (\text{Лш.5.20})$$

Наряду с задачей (УШ.Г.1-2) рассмотрим задачу $\|x - x_0\| = r$ и, $x \in M_-$, где M_- определено формулой (УШ.5.18). Решение этих задач соответственно обозначим через x_* и y_* . Пусть еще $d = \|x_0 - x_*\|$ и $d_- = \|x_0 - y_*\|$. Так как $M_- \subset M$, то $d \leq d_-$; кроме того, $x_* \in M$ и поэтому $Qx_* \in M_-$. Отсюда, в силу неравенства (Лш.5.20),

$$d \leq d_- \leq \|x_0 - Qx_*\| \leq \|x_0 - x_*\| + \|x_* - Qx_*\| = d + c\delta \|x_*\|. \quad (\text{УШ.5.21})$$

Имеем разность $x_* - y_*$. Имеем

$$\|x_* - y_*\| \leq \|x_* - Qx_*\| + \|Qx_* - y_*\|. \quad (\text{УШ.5.22})$$

Точка $Qx_* \in M_-$; по лемме Лш.2.1 $\|Qx_* - y_*\|^2 \leq \|x_0 - Qx_*\|^2 - d_-^2$. Теперь по неравенству (УШ.5.21) $\|Qx_* - y_*\|^2 \leq c\delta \|x_*\| (2d + c\delta \|x_*\|)$. Подставив это в (УШ.5.22), получим

$$\|x_* - y_*\| \leq c\delta \|x_*\| \delta + \sqrt{c\delta \|x_*\| (2d + c\delta \|x_*\|)} \sqrt{\delta}; \quad (\text{УШ.5.23})$$

если элемент x_* фиксирован, то отсюда получается более простое неравенство

$$|x_* - y_*| < C\sqrt{\delta}, \quad C = \text{const}, \quad (\text{УШ.5.24})$$

которое доказывает устойчивость решения задачи (УШ.1.1-2) по отношению к малому уменьшению числа d , если $\|q\| < d$.

Пусть теперь d заменено на $d + \delta$. Обозначим

$$M_+ = \{x \in H : \|x - q\| \leq d + \delta\}. \quad (\text{УШ.5.25})$$

Если x_* - решение задачи

$$|x - x_0| = \min, \quad x \in M_+, \quad (\text{УШ.5.26})$$

и $d_+ = |x_0 - x_*|$, то справедливы соотношения, аналогичные соотношениям (УШ.5.21) и (УШ.5.23):

$$d_+ \leq d \leq d_+ + C\delta |x_*|, \quad (\text{УШ.5.27})$$

$$|x_* - x_*| \leq C|x_*|\delta + \sqrt{C|x_*|\delta + C^2\delta^2} \sqrt{\delta}. \quad (\text{УШ.5.28})$$

В неравенстве (УШ.5.28) заменим левую часть меньшей величиной $|x_*| - |x_*|$. Теперь легко получить для $|x_*|$ квадратное неравенство

$$|x_*|^2 - 2|x_*| \frac{|x_*|(1-C\delta) + 2C\delta}{1-2C\delta - C^2\delta^2} + \frac{|x_*|^2}{1-2C\delta - C^2\delta^2} \leq 0,$$

из которого следует, что

$$|x_*| \leq (1-2C\delta - C^2\delta^2)^{-1} [(1-C\delta)|x_*| + 2C\delta + \sqrt{2C^2\delta^2|x_*|^2 + 2C\delta\delta|x_*| + 4C^2\delta^2\delta^2}].$$

Для малых δ последний радикал есть величина порядка

$$O(|x_*|\delta + \sqrt{|x_*|\delta + \delta}) \text{ и, следовательно, } |x_*| \leq |x_*| + O(\sqrt{|x_*|\delta}).$$

Подставив это в (УШ.5.28), убедимся, что при фиксированном x_*

$$|x_* - x_*| = O(\sqrt{\delta}); \quad (\text{УШ.5.29})$$

это доказывает устойчивость решения рассматриваемой задачи относительно малого уменьшения числа d , если $\|q\| < d$.

6°. Принимая по-прежнему, что $\|g\| < d$, допустим, что g заменено элементом g' , таким, что $\|g - g'\| \leq \delta$; остальные данные не изменены. Имеем новую задачу

$$\|x - x_0\| = \min, x \in M', M' = \{x \in H : \|x - g'\| \leq d\}. \quad (\text{Уш. 5.30})$$

При этом

$$\|x - g\| - \|g - g'\| \leq \|x - g'\| \leq \|x - g\| + \|g - g'\|,$$

или

$$\|x - g\| - \delta \leq \|x - g'\| \leq \|x - g\| + \delta. \quad (\text{Уш. 5.31})$$

Обозначим, как и выше,

$$M_{\pm} = \{x \in H, \|x - g\| \leq d \pm \delta\}. \quad (\text{Уш. 5.32})$$

Пусть $x \in M_-$. Тогда $\|x - g'\| \leq \|x - g\| + \delta \leq d$, т.е. $x \in M'$. Таким образом, $M_- \subset M'$. Точно так же доказывается, что $M' \subset M_+$. Окончательно, неравенство (Уш. 5.31) означает, что

$$M_- \subset M' \subset M_+. \quad (\text{Уш. 5.33})$$

При расширении множества минимум уменьшается, поэтому (смысл обозначений очевиден)

$$d_+ \leq d' \leq d_-, d_+ \leq d \leq d_-. \quad (\text{Уш. 5.34})$$

Из формул (Уш. 5.21) и (Уш. 5.27) следует, что $d_- - d_+ = O(\delta^2)$, а тогда тот же порядок имеет и разности $d_- - d'$, $d_- - d$, $d' - d_+$, $d - d_+$.

Пусть, далее, x'_* и x_* - минимизирующие элементы, соответствующие множествам M' и M_+ . Так как $M' \subset M_+$, то $x'_* \in M_+$; по лемме Уш. 2.1, $|x'_* - x_*|^2 \leq d'^2 - d_+^2$ и величина $|x'_* - x_*| = O(\sqrt{\delta})$. Точно так же $x_* \in M \subset M_+$, и $|x_* - x'_*| = O(\sqrt{\delta})$. Теперь

$$\|x_* - x'_*\| \leq \|x_* - x_*\| + \|x_* - x'_*\| = O(\sqrt{\delta}); \quad (\text{Уш. 5.35})$$

решение задачи (Уш. I.1-2) оказывается устойчивым по отношению к малым (в метрике $\|\cdot\|$) возмущениям элемента g .

7°. Пусть одновременно мало возмущены все параметры задачи (Уш. I.1-2): элемент x_0 заменен на x_1 , норма $\|\cdot\|$ - на $\|\cdot\|_1$, число d заменено на $d = d \pm \delta$ и, наконец, элемент g заменен близким,

в матрике $\| \cdot \|$, элементом Q' . По-прежнему пусть X_* означает решение невозмущенной задачи (Уш. I. I-2); через $X_{1*}, X_{2*}, X_{3*}, X_{4*}$ обозначим решения следующих возмущенных задач: X_{1*} соответствует элементу X_1 при остальных невозмущенных параметрах, X_{2*} соответствует элементу X_2 и норме $\| \cdot \|_1$ при невозмущенных Q и Q' , X_{3*} - решение, полученное при невозмущенном Q и остальных возмущенных параметрах и, наконец, X_{4*} соответствует возмущенным значениям всех четырех параметров. По-прежнему положим $d = |x_0 - x_*|$; пусть, кроме того, $d_i = |x_i - x_{i*}|$, $i = 1, 2, 3, 4$. Оценим разность $X_* - X_{4*}$. Имеем

$$|x_* - x_{4*}| \leq \sum_{i=1}^4 |x_{i-1*} - x_{i*}|, \quad x_{0*} = x_*;$$

в силу неравенств (Уш. 5.2-3)

$$|x_* - x_{4*}| \leq \delta + \sum_{i=2}^4 |x_{i-1*} - x_{i*}|. \quad (\text{Уш. 5.2})$$

Как видно из формул (Уш. 5.6), (Уш. 5.29) и (Уш. 5.35), оценки норм в (Уш. 5.36) определяются через оценки величин d_i в зависимости от d . Найдем эти последние оценки. Очевидно,

$$d_1 \leq |x_1 - x_0| + |x_0 - x_*| + |x_* - x_{1*}| \leq d + 2\delta; \quad (\text{Уш. 5.2})$$

по формулам (Уш. 5.3), (Уш. 5.8) и (У. 5.37)

$$\begin{aligned} d_2 &= |x_2 - x_{2*}| \leq (1-\delta) |x_1 - x_{2*}| \leq (1-\delta) (|x_1 - x_{1*}| + |x_{1*} - x_{2*}|) \leq \\ &\leq (1+\delta)(1-\delta) (d+2\delta + C_5(d+2\delta)\sqrt{\delta}) = d + O(\sqrt{\delta}). \end{aligned} \quad (\text{Уш. 5.2})$$

Оценим величину d_3 . По формулам (Уш. 5.27), в случае $d' = d + \delta$ будет $d_3 = d_2 + O(\delta) = d + O(\sqrt{\delta})$; то же соотношение верно и для $d' = d - \delta$ - это следует из формул (Уш. 5.21). Теперь из (Уш. 5.34) вытекает, что $d_4 = a_3 + O(\sqrt{\delta}) = d + O(\sqrt{\delta})$.

Даже проверить, что все слагаемые в (Уш. 5.36) суть величины порядка $O(\sqrt{\delta})$ - это означает, что решение X_* устойчиво относительно малых возмущений всех параметров, перечисленных в начале настоящего пункта.

Аналогичными рассуждениями доказывается устойчивость решения задачи (Уш. 5.1) относительно одномерных малых возмущений

элемента X_0 , нормы $|\cdot|$ и множества M , если последнее возмущению соответствует условиям п. 4^а.

§ 6. Оценка погрешности искажений

1^о. Пусть $a(x, y)$ симметричная билинейная форма, определенная на некотором линейном множестве \mathcal{L} , и соответствующая ей квадратичная форма $a(x, x) = a(x, x)$ положительна, - это значит, что $a(x, x) \geq 0$ и $a(x, x) = 0 \Rightarrow x = 0$. Введя скалярное произведение $[x, y] = a(x, y)$ и норму $|x| = \sqrt{a(x, x)}$, мы превратим множество \mathcal{L} в предгильбертово пространство; замкнув по лемме в указанной норме, получим гильбертово пространство, которое обозначим через H . Примем, что H сепарабельно; для упрощения последующих записей примем еще, что оно вещественно.

2^о. Пусть $f(x)$ - линейный ограниченный в H функционал, определенный на всем пространстве, и $M \subset H$ - непустое выпуклое множество, замкнутое в норме H . Задача

$$a(x) - 2f(x) = |x|^2 - 2f(x) = \min, \quad x \in M. \quad (\text{VII.6.1})$$

имеет одно и только одно решение: она равносильна задаче

$$|x - X_0| = \min, \quad x \in M, \quad (\text{VII.6.2})$$

где X_0 - элемент пространства H определенный тождеством

$$\forall x \in H, \quad f(x) = [x, X_0]; \quad (\text{VII.6.3})$$

существования и единственности элемента X_0 вытекают из известной теоремы Ф. Риса.

3^о. Приближенное решение задачи (VII.6.1) можно строить, например, по методу Рунда (в книге Гловинского, Дворжа и Трельмьера [1] вместо термина "метод Рунда" используется термин "внутренние методы"). Зададим полную в H последовательность конечномерных подпространств $\{H_n\}$, и пусть $(y_{1n}, y_{2n}, \dots, y_{m_n})$, $N = N(n)$ - базис подпространства H_n . Задачу (VII.6.1) заменим следующей

$$|x^{(n)}|^2 - 2f(x^{(n)}) = \min, \quad x^{(n)} \in H_n \cap M, \quad (\text{VII.6.4})$$

т.е., что то же,

$$|x^{(n)} - x_0| = \min, \quad x^{(n)} \in H_n \cap M. \quad (\text{Уш. 6.})$$

Если пересечение $H_n \cap M$ не пусто, то задача (Уш. 6.4) имеет одно и только одно решение.

4°. Пусть

$$x = \sum_{k=1}^{\mathcal{N}} a_k^{(n)} y_{kn}, \quad (\text{Уш. 6.})$$

тогда

$$|x^{(n)}|^2 = \sum_{j,k=1}^{\mathcal{N}} [y_{jn}, y_{kn}] a_j^{(n)} a_k^{(n)}. \quad (\text{Уш. 6.})$$

Скалярные произведения $[y_{jn}, y_{kn}]$ вычисляются с некоторыми погрешностями, которые можно считать сколь угодно малыми. Обозначим погрешность вычисления величинами $[y_{jn}, y_{kn}]$ через $y_{jk}^{(n)}$ и положим $[y_{jn}, y_{kn}] + y_{jk}^{(n)} = d_{jk}^{(n)}$; числа $d_{jk}^{(n)}$ нам известны. Обозначим соответственно через M_n и Γ_n матрицы элементов $[y_{jn}, y_{kn}]$ и $y_{jk}^{(n)}$; $j, k = 1, 2, \dots, \mathcal{N}$, и через $\|M_n\|, \|\Gamma_n\|$ нормы этих матриц как операторов в $R_{\mathcal{N}}$.

Матрица M_n симметрична, и мы вправе считать, что матрица Γ_n также симметрична. Можно считать также, что матрица $M_n + \Gamma_n$ положительно определенная - это будет всегда, когда Γ_n достаточно мала по отношению к M_n .

Если $x^{(n)}$ - произвольный элемент из H_n ,

$$x^{(n)} = \sum_{k=1}^{\mathcal{N}} C_k^{(n)} y_{kn},$$

то по формуле (Уш. 6.7)

$$|x^{(n)}|^2 = (M_n C^{(n)}, C^{(n)}),$$

где $C^{(n)} = (C_1^{(n)}, C_2^{(n)}, \dots, C_{\mathcal{N}}^{(n)})$. Искожаемая матрица $\tilde{M}_n := M_n + \Gamma_n$ порождает в H_n новую норму

$$|x^{(n)}|_1^2 := (\tilde{M}_n C^{(n)}, C^{(n)}) = |x^{(n)}|^2 + (\Gamma_n C^{(n)}, C^{(n)}).$$

Примем, что матрица Γ_n мала по сравнению с M_n в следующем смысле:

$$\forall C^{(n)} \in R_{\mathcal{N}}, \quad \frac{|(\Gamma_n C^{(n)}, C^{(n)})|}{|(M_n C^{(n)}, C^{(n)})|} \leq \delta_1, \quad (\text{Уш.})$$

где δ_1 - малое число. Заметим, что условие (Уш.6.8) выполнено, если

$$\|\Gamma_n\| \cdot \|M_n^{-1}\| \leq \delta_1. \quad (\text{Уш.6.9})$$

Действительно, обозначим $M_n^{1/2} C^{(n)} = \rho^{(n)}$, тогда отношение (Уш.6.8) принимает вид

$$\frac{|\langle \Gamma_n M_n^{-1/2} \rho^{(n)}, M_n^{-1/2} \rho^{(n)} \rangle|}{\|\rho^{(n)}\|^2} \leq \|\Gamma_n\| \cdot \|M_n^{-1/2}\|^2 = \|\Gamma_n\| \cdot \|M_n^{-1}\|.$$

Если соблюдено условие (Уш.6.8), то

$$(1-\delta)|x^{(n)}| \leq |x^{(n)}|_1 \leq (1+\delta)|x^{(n)}|, \quad (\text{Уш.6.10})$$

где $\delta = \max(\sqrt{1+\delta_1}-1, 1-\sqrt{1-\delta_1}) = \delta_1/(1+\sqrt{1-\delta_1})$

ξ^0 . Положим $X_0 = X_{on} + X_m$, где $X_n \in H_n$ и $X_m \perp H_n$. Задача (Уш.6.5) равносильна такой задаче:

$$|x^{(n)} - x_{on}| = \min, \quad x \in H_n \cap M; \quad (\text{Уш.6.11})$$

элемент X_{on} можно построить как проекцию X_0 на H_n , т.е., как решение экстремальной задачи

$$|x_0 - \sum_{k=1}^N b_k y_{kn}|^2 = \min$$

что приводит к системе уравнений

$$\sum_{k=1}^N b_k [y_{jn}, y_{jn}] = [x_0, y_{jn}] = f(y_{jn}), \quad 1 \leq j \leq N. \quad (\text{Уш.6.12})$$

Эту систему можно записать в виде

$$[x_{on}, y_{jn}] = f(y_{jn}), \quad 1 \leq j \leq N;$$

таким образом, X_{on} есть элемент, даваемый теоремой Э.Риса, тогда функционал f считать на H_n .

При оставлении системы (Уш.6.12) будем допускать погрешности: вместо $[y_{kn}, y_{jn}]$ пишем коэффициенты $[y_{kn}, y_{jn}] + \delta_{kj}^{(n)} = d_{kj}$; кроме того, числа $f(y_{jn})$ также будем считать с погрешностями: пусть последние члены $\delta_j^{(n)}$. Вместо системы (Уш.6.12) мы в следующем разделе получим искаженную систему

$$M_n \tilde{b}^{(n)} = f^{(n)} + \delta^{(n)}. \quad (\text{Уш.6.13})$$

Здесь $\tilde{b}^{(n)}$ - искаженный вектор коэффициентов,

$$\tilde{b}^{(n)} = (\tilde{b}_1^{(n)}, \tilde{b}_2^{(n)}, \dots, \tilde{b}_M^{(n)}); \quad (\text{Уш. 6.})$$

$$\delta^{(n)} = (\delta_1^{(n)}, \delta_2^{(n)}, \dots, \delta_M^{(n)}).$$

Допустим, что линейный вычислительный процесс определен элемента X_0 по методу Рунге устойчив в последовательности (R_M) ; напомним, что для этого необходимо и достаточно, чтобы координатная система была сильно минимальна в H . Тогда существуют такие положительные постоянные ρ, q, τ , что если $\|\Gamma_n\|$ то

$$\|y_{on} - x_{on}\| < \rho \|\Gamma_n\| + q \|\delta^{(n)}\|;$$

здесь y_{on} - искаженное значение элемента x_{on} . Пусть $\|\Gamma_n\|$ и $\|\delta^{(n)}\|$ достаточно малы. Тогда $\|y_{on} - x_{on}\| < \delta$, где δ - малое число не зависящее от n . Одновременно, в силу необходимых условий устойчивости, нормы $\|\Gamma_n\|$ равномерно ограничены, и неравенства (Уш. 6.9-10) также выполнены. Таким образом, искажение сводится к тому, что при каждом n возмущаются элемент x_{on} и норма в пространстве H_n , и эти возмущения малы равномерно относительно n . По доказанному в § 5, искажение элемента x_{on} приводит к искажению решения на $O(\|\Gamma_n\| + \|\delta^{(n)}\|)$, а искажение нормы - к искажению решения на $O(\sqrt{\|\Gamma_n\|} \cdot \|\delta^{(n)}\|) = O(\sqrt{\|\Gamma_n\|})$. Следовательно, если $y_*^{(n)}$ - искаженное решение n -й приближенной задачи

$$\|y_*^{(n)} - x_*^{(n)}\| = O(\sqrt{\|\Gamma_n\|} + \|\delta^{(n)}\|). \quad (\text{Уш. 6.})$$

Как показывает неравенство (Уш. 6.15), с целью избежать больших искажений следует вычислять матрицу Рунге точнее, чем свободные члены системы Рунге - так, чтобы $\|\Gamma_n\| = O(\|\delta^{(n)}\|^2)$.

Можно исследовать устойчивость процесса вычисления коэффициентов $a_k^{(n)}$ не в последовательности (R_M, R_M) , а в последовательности (\bar{X}_M, \bar{Y}_M) (см. формулу Ш.3.4). Сформулируем результат упомянутый процесс устойчив и при достаточно малых нормах

верна оценка

$$\|y_*^{(n)} - x_*^{(n)}\|_{\bar{X}_M \rightarrow \bar{Y}_M} = O(\sqrt{\|\Gamma_n\|_{\bar{X}_M \rightarrow \bar{Y}_M}} + \|\delta^{(n)}\|_{\bar{Y}_M}). \quad (\text{Уш. 6.})$$

Эта оценка пригодна, в частности, для метода конечных элементов.

Дополнение к главе VIII

В настоящем дополнен и, написанном В.Б.Тихтиным излагаются результаты некоторых работ последних лет, посвященных оценкам аппроксимации для решений односторонних вариационных задач, отличных, как правило, от простейшей задачи, изученной в данной главе. Приближенные решения в этих работах обычно строятся по МКЭ.

В Дополнении приняты несколько иные термины и обозначения, чем в основном тексте гл. VIII. Соболевские пространства $W_2^{(n)}$; $W_2^{(n)}$ обозначаются соответственно через H_0^n и H_0^n . Множество ограничений обозначается через \mathcal{K} . Если ограничительное условие имеет вид $u(x) \leq \chi(x)$ или $u(x) \geq \chi(x)$, то множество ограничений называется препятствием; различаются препятствия в области Ω на границе, в зависимости от того, где выполняется упомянутое неравенство. Далее, O означает точку евклидова пространства; функции, заданные в области евклидова пространства, чаще всего обозначаются буквами u и U . Задачи, как правило, формулируются в терминах вариационных неравенств.

I⁰. В статье Балка [1] рассматривается плоская задача о препятствии в области. Пусть $\Omega \subset R_2$ — выпуклая область с границей класса $C^{(2)}$, и (по повторяющимся индексам производится суммирование от 1 до 2)

$$a(u, v) = \int_{\Omega} (a_{ij} \frac{\partial u}{\partial x_i} \frac{\partial v}{\partial x_j} + cuv) dx,$$

где

$$a_{ij}(x) \in C^{(1)}(\bar{\Omega}), \quad a_{ij} = a_{ji}, \quad c(x) \in L_{\infty}(\Omega);$$

$$a_{ij} \xi_i \xi_j \geq \alpha |\xi|^2, \quad \alpha > 0, \quad c(x) \geq \lambda > 0.$$

Для $f \in L_1(\Omega)$ рассматривается задача об отыскании такой функции $u \in \mathcal{K}$, что

$$a(u, v-u) \geq \int_{\Omega} f(v-u) dx, \quad \forall v \in \mathcal{K}, \quad (\text{УШ. I})$$

где

$$\mathcal{K} = \{ v \in H_0^1(\Omega) : v(x) \geq \chi(x) \text{ почти всюду в } \Omega \};$$

$\chi \in H^1(\Omega)$ - заданная функция, $\chi|_{\partial\Omega} = 0$.

При перечисленных предположениях решение задачи принадлежит к классу $H^1(\Omega)$ (см. Брезко и Стампакья [1]).

Для приближенного решения задачи (УШ.1) применяется МКЭ с кусочно линейными координатными функциями. Пусть Ω_h - вписанный в Ω многоугольник, длины сторон которого не превосходят h . Этот многоугольник разбивается на треугольники со сторонами $\leq h$. Через V_h обозначим пространство функций, заданных на Ω_h , непрерывных, кусочно линейных и равных нулю на $\Omega \setminus \Omega_h$. Для любой функции $v \in C(\bar{\Omega})$ равной нулю на $\partial\Omega$, обозначим через v_I ее интерполянт, т.е., функцию из V_h , совпадающую с v в вершинах треугольника. Приближенное решение задачи (УШ.1) определяется как решение задачи

$$u_h \in \mathcal{K}_h; a(u_h, v_h - u_h) \geq \int_{\Omega} f(v_h - u_h) dx; \forall v_h \in \mathcal{K}_h, \quad (\text{УШ.2})$$

где $\mathcal{K}_h = \{v_h \in V_h : v_h \geq \chi_I \text{ в узлах триангуляции}\}$. При некоторых "условиях регулярности" описанного здесь варианта МКЭ устанавливается оценка

$$\|u - u_h\|_{H^1(\Omega)} \leq ch \|u\|_{H^2(\Omega)}. \quad (\text{УШ.3})$$

В работе [2] так обобщили результаты работы [1] на случай невыпуклой области $\Omega \subset R^2$. Предположим, что $\partial\Omega \in C^{(2)}$ и имеет лишь конечное число точек, в которых кривизна меняет знак. Тогда оценка (УШ.3) по-прежнему имеет место.

2°. В статье Скарпиччи и Вивальди [1] рассмотрена плоская задача с препятствием на границе. Пусть $\Omega \in R^2$ - выпуклая ограниченная область с границей класса $C^{(2)}$, $a(u, v)$ билинейная форма, такая же, как в п.1°, $f \in L_2(\Omega)$, $\psi \in H^2(\Omega)$. Рассмотрим задачу (УШ.1), но с другим препятствием, именно,

$$\mathcal{K} = \{v \in H^1(\Omega) : v \geq \psi \text{ на } \partial\Omega\}. \quad (\text{УШ.4})$$

Известно (см. Брезко [1]), что при указанных условиях решение задачи (УШ.1), (УШ.4) $u \in H^2(\Omega)$.

Для приближенного решения задачи применен метод, близкий к методу п.1°. Выделим в Ω многоугольник Ω_h , длины сторон которого не превосходят h . Разобьем Ω_h на треугольники. Обозначим через $I_0 = \{1, 2, \dots, N_0\}$ indexes внутренних узлов, через $I_1 =$

$= \{N_0 + 1, \dots, N\}$ - индексом граничных узлов; сами узлы обозначим через X_i ; положим еще $I = I_0 \cup I_1$. Далее, пусть Γ_i - часть $\partial\Omega$ между узлами X_i и X_{i+1} , Σ_i - часть Ω , ограниченная отрезком $[X_i, X_{i+1}]$ и кривой Γ_i .

Если T_i - треугольник со стороной $[X_i, X_{i+1}]$, $i \in I_1$, то объединение $T_i \cup \Sigma_i$ рассматривается как криволинейный элемент.

Через $y_i^h(x)$, $x \in \Omega$, $i \in I$, обозначена непрерывная в Ω функция, линейная на каждом треугольнике или криволинейном элементе, равная единице в узле X_i и нулю во всех остальных узлах.

Пусть

$$H_h^1(\Omega) = \{v_h : v_h = \sum_{i \in I} v_i^h y_i^h(x), v_i^h \in P_1, \forall i \in I\}. \quad (VII.5)$$

Приближенная задача имеет вид (VII.2), но \mathcal{K}_h определяется формулой

$$\mathcal{K}_h = \{v_h \in H_h^1 : v_h(x_i) \geq \psi(x_i), \forall i \in I\}. \quad (VII.6)$$

При условии регулярности использованного МКЭ верна оценка

$$\|u - u_h\|_{H^1(\Omega)} \leq Ch^{3/4} \|u\|_{H^2(\Omega)}. \quad (VII.7)$$

3°. В статье Бреши, Хагера и Ратцлера [1] рассматривается конечно-элементная аппроксимация задачи с препятствием в области или на границе для некоторых нелинейных форм.

Задача с препятствием в области ставится в форме

$$u \in \mathcal{K};$$

$$\int_{\Omega} [\nabla u \cdot \nabla(v-u) + u(v-u)] dx \geq \int_{\Omega} f(v-u) dx; \quad \forall v \in \mathcal{K}; \quad (VII.8)$$

$$\mathcal{K} = \{v \in H^1(\Omega) : v \geq \psi \text{ в } \Omega, v|_{\partial\Omega} = g\}.$$

Здесь Ω - выпуклая область класса $C^{(1,1)}$; $g, \psi \in H^1(\Omega)$, $\psi|_{\partial\Omega} \leq g$, $f \in L_2(\Omega)$.

Пусть Ω_h - вписанный в Ω многоугольник со сторонами, по длине не превосходящими h , и пусть этот многоугольник регулярным образом триангулирован. Для каждого треугольника триангуляции T_i , прилегающего к $\partial\Omega_h$, построим треугольник T_i' зеркальное отражение треугольника T_i в его стороне, принадлежащей

$\partial\Omega_h$. Пусть $R_i = UT_i'$. Будем считать h настолько малым, что $(\Omega_h \cup R_h) \subset \Omega$. Обозначим через Ω некоторую область, содержащую $\Omega_h \cup R_h$ при любом достаточно малом h и через V_h - множество непрерывных кусочно линейных на $\Omega_h \cup R_h$ функций. Продолжим решение $u \in H^1(\Omega)$ задачи (УШ.7) до функции $\tilde{u} \in H^1(\Omega)$. Пусть \tilde{u}_I, Ψ_I - интерполанты функций \tilde{u} и Ψ соответственно на $\Omega_h \cup R_h$. Множество \mathcal{K}_h определяется следующим образом:

$$\mathcal{K}_h = \{v \in V_h : v_h \geq \Psi_I \text{ в узлах на } \Omega_h, v_h = \tilde{u}_I \text{ на } \Omega \setminus \Omega_h\}.$$

Приближенное решение u_h определяется согласно (УШ.7) с заменой \mathcal{K} на \mathcal{K}_h . Устанавливается оптимальная по порядку оценка погрешности аппроксимации:

$$\|u - u_h\|_{H^1(\Omega)} = O(h). \quad (\text{УШ.9})$$

Наряду с кусочно линейными координатными функциями для приближенного решения задачи (УШ.7) применяются также кусочно квадратичные функции. В этом случае, в предположении, что точное решение $u \in H^{3/2-\epsilon}(\Omega)$, устанавливается оценка

$$\|u - u_h\|_{H^1(\Omega)} = O(h^{3/2-\epsilon}). \quad (\text{УШ.10})$$

Задача о вращении на границе также ставится в форме (УШ.8), но \mathcal{K} определяется формулой $\mathcal{K} = \{v \in H(\Omega), v|_{\partial\Omega} \geq g\}$.

Пространство V_h определяется следующим образом. Если $v_h \in V_h$, то v_h непрерывна и кусочно линейна на Ω_h ; при этом, если точка $Q \in \Omega \setminus \Omega_h$, то $v_h(Q) = v_h(P)$, где P - проекция Q на $\partial\Omega_h$.

Положим в данном случае

$$\mathcal{K}_h = \{v_h \in V_h, v_h \geq g \text{ в узлах на } \partial\Omega\}.$$

Для любых вале решение строится как решение задачи (УШ.8), в которой \mathcal{K} заменено на \mathcal{K}_h . В предположении, что точное решение

$u \in H^1(\Omega)$, а $u, g \in W^1$, вблизи $\partial\Omega$ и свободная граница свободной конечного числа точек, устанавливается оптимальная по порядку оценка погрешности аппроксимации:

$$\|u - u_h\|_{H^1(\Omega)} = O(h). \quad (\text{УШ.11})$$

Таким образом, при более сильных ограничениях на решение получается более высокая по порядку оценки погрешности сравнительно с работой Скардина и Живальди [7].

4°. В статье ниже [1] рассматривается задача о минимуме функционала

$$I(u) = \int_{\Omega} |\nabla u|^2 dx - 2 \int_{\Omega} f u dx \quad (\text{Уш. 12})$$

на множестве $\mathcal{K} = \{u : u \in H_0^1(\Omega), u \leq \chi \text{ почти всюду в } \Omega\}$.

Здесь $\Omega \in R_2$ - область с достаточно гладкой границей, $\chi \in W_{\infty}^{(2)}(\Omega)$

и $\chi, f \in L_2(\Omega)$ - заданные функции.

Рассматривается регулярная триангуляция с параметром h , области Ω и строится пространство S_h непрерывных кусочно-линейных функций, равных нулю на $\partial\Omega$. Пусть $\mathcal{K}_h = \mathcal{K} \cap S_h$. Для задачи (Уш. 12) приближенная задача состоит в минимизации функционала $I(u)$ на множестве \mathcal{K}_h . Пусть u и u_h соответственно точное и приближенное решение. Если $u \in W_{\infty}^{(2)}(\Omega)$, то устанавливается следующая оценка погрешности аппроксимации:

$$\|u - u_h\|_{W_{\infty}(\Omega)} \leq Ch^2 |\ln h| \cdot \|u\|_{W_{\infty}^{(2)}(\Omega)}, \quad (\text{Уш. 13})$$

где $C > 0$ не зависит от h .

5°. В статье Лительмана [1] рассматривается задача о максимуме в области для некоторого неквадратичного функционала. Именно, пусть $\Omega \in R_n$, $\mathcal{K} = \{u \in W_2^{(1)}(\Omega) : u(x) \geq \psi(x) \text{ в } \Omega,$

$u|_{\partial\Omega} = f\}$, где ψ и f заданные функции, причем $\psi|_{\partial\Omega} \leq f$. Рассматривается задача о минимизации функционала

$$\int_{\Omega} F(\nabla u) dx \quad (\text{Уш. 14})$$

на множестве \mathcal{K} . В интеграле (Уш. 14) $F \in C^{(2)}(R_n)$.

Позначим $\alpha(p) = \nabla F(p) = (\alpha_1(p), \alpha_2(p), \dots, \alpha_n(p))$, $p \in R_n$

и допустим, что $(\alpha(p) - \alpha(q)) \cdot (p - q) \geq c \max(|p|, |q|) |p - q|$

$\forall p, q \in R_n$, $c > 0$ - евклидова норма в R_n . Тогда задача

(Уш. 14) равносильна задаче

$$u \in \mathcal{K};$$

$$a(u, v-u) = \int_{\Omega} \alpha_i(\nabla u) \frac{\partial}{\partial x_i} (v-u) dx > 0, \forall v \in \mathcal{K};$$

эта последняя аппроксимируется по МКЭ с кусочно линейными и квадратными функциями. Пусть $\partial\Omega \in C^{(2)}$. Обозначим через T_1, T_2, \dots, T_m симплексы разбиения и через Ω_h - объединение $\bigcup_{i=1}^m T_i$. h обозначен параметр, которого не превосходят длины ребер элементов T_i . Предполагается, что внутренности симплексов T_i не пересекаются, что каждая грань любого симплекса либо есть грань другого симплекса, либо ее вершины лежат на $\partial\Omega$. V_h обозначено пространство функций, непрерывных на Ω_h и линейных на каждом симплексе T_i . Обозначим еще через $\mathcal{N} = \mathcal{N}_h$ множество узлов в Ω_h . Пусть

$$\mathcal{K}_h = \{v_h \in S_h : v_h(x) \geq \psi(x), x \in \Omega \cap \mathcal{N}; v_h(x) = f(x), x \in \partial\Omega \cap \mathcal{N}\}$$

Приближенное решение определяется как решение следующей задачи

$$u_h \in \mathcal{K}_h,$$

$$\int_{\Omega} \alpha_i(\nabla u_h) \frac{\partial}{\partial x_i} (v_h - u_h) dx \geq 0, \forall v_h \in \mathcal{K}_h;$$

в предположении, что $u \in H^1(\Omega) \cap W_{\infty}^{(1)}(\Omega)$, устанавливается оценка погрешности

$$\|u - u_h\|_{H^1(\Omega_h)} = O(h). \quad (1)$$

6°. В работе Л. С. Далецкого и Марро [1] рассмотрена следующая задача. Пусть Ω - выпуклая область в R_n и $p > 1$. Требуется найти функцию $u \in \dot{W}_p^{(1)}(\Omega)$, такую, что

$$J(u) = \inf_{v \in \dot{W}_p^{(1)}(\Omega)} J(v),$$

$$J(v) = \frac{1}{p} \int_{\Omega} |\nabla u|^p dx - \int_{\Omega} f v dx; f \in W_p^{(1)}(\Omega), \frac{1}{p} + \frac{1}{p'} = 1.$$

Пусть V_h - семейство конечномерных подпространств в $\dot{W}_p^{(1)}$. Приближенное решение задачи (7.16) определим как решение задачи

$$J(u_h) = \inf_{v_h \in V_h} J(v_h).$$

Задачи (УШ.16) и (УШ.17) имеют единственные решения u и u_h соответственно (см., например, Якоби [2]).

Пусть $E_h(u) = \inf_{v_h \in V_h} \|v_h - u\|_{W_p^{(1)}(\Omega)}$. Устанавливаются следующие оценки погрешности аппроксимации:

$$\|u_h - u\|_{W_p^{(1)}(\Omega)} \leq C [E_h(u)]^{1/(p-1)}, \quad p \geq 2; \quad (\text{УШ.18})$$

$$\|u_h - u\|_{W_p^{(1)}(\Omega)} \leq C [E_h(u)]^{1/(3-p)}, \quad 1 < p < 2,$$

где C не зависит от h .

Оценки (УШ.18) улучшены в работе Тюрина [2]:

$$\|u_h - u\|_{W_p^{(1)}(\Omega)} \leq C [E_h(u)]^{2/p}, \quad p \geq 2; \quad (\text{УШ.19})$$

$$\|u_h - u\|_{W_p^{(1)}(\Omega)} \leq C [E_h(u)]^{p/2}, \quad 1 < p < 2.$$

7°. Пусть Ω - единичный квадрат на плоскости: $\Omega = [0, 1] \times [0, 1]$; на Ω заданы функции u_0, v_0 , причем $u_0 \leq v_0$ почти всюду в Ω , $u_0|_{\partial\Omega} \leq 0 \leq v_0|_{\partial\Omega}$, $p > 2$

$$J(u) = \frac{1}{p} \sum_{i=1}^2 \int_{\Omega} \left| \frac{\partial u}{\partial x_i} \right|^p dx - \int_{\Omega} f u dx. \quad (\text{УШ.20})$$

В работе Скарпини [1] рассматривается, в частности, проблема оценки приближенного решения для задачи минимизации функционала (УШ.20) на множестве:

$$K = \left\{ v \in W_p^{(1)}(\Omega) : u_0(x) \leq v(x) \leq v_0(x) \text{ почти всюду в } \Omega \right\}.$$

Пусть u - точное решение. Если $f \in L_p(\Omega)$; $u_0, v_0 \in W_p^{(1+\sigma)}(\Omega)$, $\sigma = 1/(p-1)$, то $u \in W_p^{(1+\sigma)}(\Omega)$ (ср. Эль Колли [1], Моско [1]).

Пусть u_h - приближенное решение, полученное по МКЭ при кусочно-линейных координатных функциях. Докажем, что

$$\|u_0 - u_h\|_{W_p^{(1)}(\Omega)} \leq C h^{\sigma^2}, \quad (\text{УШ.21})$$

где C не зависит от h .

8°. Статья Эдмонта [1] посвящена следующей задаче. Пусть $\Omega \subset R_2$ - выпуклый многоугольник. Требуется найти такую функцию $u \in H^1(\Omega)$, что

$$J(u) := \alpha \int_{\Omega} u^2 dx + \beta \int_{\Omega} |\nabla u|^2 + j(u) - \int_{\Omega} f u dx = \min, \quad (\text{УШ.22})$$

или, что то же,

$$\alpha(u, v-u) + \beta a(u, v-u) + j(u) \geq j(v-u);$$

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v dx, \quad (u, v) = \int_{\Omega} uv dx, \quad f \in L_2(\Omega). \quad (\text{УШ.23})$$

Здесь α, β - положительные постоянные,

$$j(v) = \int_{\Omega} \Phi(v(x)) dx,$$

а функция Φ определяется формулой

$$\Phi(\lambda) = \frac{1}{2} d_2 (\lambda - u_0)_+^2 + \frac{1}{2} d_1 (\lambda - u_0)_-^2 + L(\lambda - u_0)_+,$$

где u_0, d_1, d_2, L - положительные постоянные, $\lambda_+ = \max(\lambda, 0)$, $\lambda_- = \min(\lambda, 0)$, $\lambda \in R_1$.

Существование и единственность решения задачи (УШ.22) следуют из общих теорем (см. Lions [1]).

Задача (УШ.22) аппроксимируется по МКЭ с кусочно линейными координатными функциями. Область Ω покрывается семейством T^h треугольников, h - наибольшая из сторон треугольников. Пусть V_h - пространство непрерывных кусочно линейных функций, равных нулю на $\partial\Omega$. Через $v_I \in V_h$ обозначается интерполант функции v по "злам" триангуляции. Обозначим еще

$$j_h(v) = \int_{\Omega} \{\Phi(v)\}_I dx = \frac{1}{6} \sum_{T \in T^h} A(T) \sum_{i=1}^3 \Phi(v_i^T),$$

где $A(T)$ есть площадь T -го треугольника, а v_i^T есть значение функции v в i -й вершине треугольника T . Определяется дискретное скалярное произведение

$$(w, v)_h = \sum_{T \in T^h} \frac{1}{3} A(T) \sum_{i=1}^3 w_i^T v_i^T.$$

Приближенная задача ставится так: найти функцию $u_h \in V_h$, удовлетворяющую дискретному вариационному равенству

$$\alpha(u_h, v_h - u_h)_h + \beta a(u_h, v_h - u_h) + j_h(v_h) - j_h(u_h) \geq (f, v_h - u_h); \quad (\text{УШ.24})$$

$$\forall v_h \in V_h.$$

Решение задачи (УШ.24) существует и единственно. При некоторых предположениях о свойствах точного решения устанавливается следующая оценка погрешности аппроксимации:

$$\|u - u_h\|_{H^1(\Omega)} \leq Ch, \quad (\text{УШ. 5})$$

где C не зависит от h .

ГЛАВА IX

Упруго-пластическое состояние по Сен-Венану - Мизесу и Хаару - Карману

Как хорошо известно, впервые математическая теория пластического состояния была предложена Сен-Венаном в 1870 г. Эта теория давала описание пластического состояния в условиях плоской деформации; на трехмерный случай ее распространил М.леви в том же 1870 г. В 1913 г. Мизес предложил описывать трехмерное пластическое состояние условием, близким к условию Сен-Венана и в то же время значительно более простым. Как выяснилось потом, условие Мизеса и это простое механическое трактовка условия Сен-Венана для плоской деформации уравнения Сен-Венана и Мизеса совпадают.

В 1909 г. Хаар и Карман опубликовали вариационный принцип, позволяющий описывать напряженное состояние в упруго-пластической среде, т.е. в среде, часть которой находится в упругом состоянии, а остальная часть - в пластическом. Принцип Хаара - Кармана основывается на условии пластичности Сен-Венана - Леви, но аналогичный принцип очень легко сформулировать, исходя из условия Мизеса.

Вследствие выяснилось, что явление пластического состояния значительно сложнее и многообразнее, и что уравнения Сен-Венана - Леви и Мизеса описывают только так называемое пластическое течение. В этом плане следует прежде всего отметить работы Хенки, Франдтля, Прагера, Холда, Хилла, Ирришина, Качанова и ряда других авторов.

В настоящей главе мы рассмотрим некоторые из задач, к которым приводит принцип Хаара - Кармана. Как мы увидим ниже, эти задачи входят под схему (Ш. I-2). I мы для облегчения отскажем формулируем уравнения Сен-Венана - Леви и Мизеса и вариационный принцип Хаара - Кармана. В § 2 рассматривается задача упруго-пластического кручения цилиндрических стержней, в § 3 исследуется плоская деформация, в § 4 - трехмерная задача упруго-пластического состояния. Мы сосредоточимся в этих параграфах исследовать устойчивость точных решений перечисленных задач с помощью малых изгибов, в них у оценки погрешности аппроксимации и скажем.

§ I. Уравнения Сен-Венана - Леви и Мизеса. Вариационный принцип Хаара - Кармана.

1°. Тензор напряжений будем обозначать так:

$$T = \begin{pmatrix} \tau_{xx} & \tau_{xy} & \tau_{xz} \\ \tau_{xy} & \tau_{yy} & \tau_{yz} \\ \tau_{xz} & \tau_{yz} & \tau_{zz} \end{pmatrix}. \quad (\text{IX.1.1})$$

Если $\vec{K} = (X, Y, Z)$ есть вектор с единичной длиной, то в состоянии равновесия тензор напряжений удовлетворяет уравнению равновесия

$$\text{div } T \vec{K} = 0 \quad (\text{IX.1.2})$$

или, в более подробной записи

$$\begin{aligned} \tau_{xx/x} + \tau_{xy/y} + \tau_{xz/z} + \nu &= 0, \\ \tau_{xy/x} + \tau_{yy/y} + \tau_{yz/z} + \mu &= 0, \\ \tau_{xz/x} + \tau_{yz/y} + \tau_{zz/z} + \lambda &= 0, \end{aligned} \quad (\text{IX.1.3})$$

индекс за косой черточкой означает дифференцирование по соответствующему направлению.

Вектор смещений будем обозначать через $U = (U_x, U_y, U_z)$, вектор скорости смещений — через V , так что $V = \dot{U} =$

$= (U_{x,t}, U_{y,t}, U_{z,t})$. Если считать деформации и их скорости малыми, тогда тензор скоростей деформации можно записать так:

$$\frac{1}{2} \begin{pmatrix} 2U_{xx} & U_{xy} + U_{yx} & U_{xz} + U_{zx} \\ U_{xy} + U_{yx} & 2U_{yy} & U_{yz} + U_{zy} \\ U_{xz} + U_{zx} & U_{yz} + U_{zy} & 2U_{zz} \end{pmatrix}. \quad (\text{IX.1.4})$$

Главные нормальные напряжения обозначим через $\tau_1 \geq \tau_2 \geq \tau_3$; они суть корни уравнения

$$\begin{vmatrix} \tau_{xx} - \tau & \tau_{xy} & \tau_{xz} \\ \tau_{xy} & \tau_{yy} - \tau & \tau_{yz} \\ \tau_{xz} & \tau_{yz} & \tau_{zz} - \tau \end{vmatrix} = 0. \quad (\text{IX.1.5})$$

Наибольшее касательное напряжение равно

$$\tau_m = \frac{\tau_1 - \tau_2}{2}. \quad (\text{IX.1.6})$$

2°. Сен-Венан [1] сформулировал условия пластичности и написал уравнения пластического состояния для случая плоской деформации; в этом случае напряжения не зависят от одной из координат, скажем от z , и $U_z = \tau_{xz} = \tau_{yz} = 0$. Одно из главных нормальных напряжений равно τ_{xz} , а остальные два равны

$$\frac{1}{2} [\tau_{xx} + \tau_{yy} \pm \sqrt{(\tau_{xx} - \tau_{yy})^2 + 4\tau_{xy}^2}]. \quad (\text{IX.1.7})$$

От работы на опыте Треска, Сен-Венан формулирует следующие гипотезы пластического состояния:

I. Наибольшее касательное напряжение в пластическом состоянии есть величина постоянная, характерная для данной среды. Сен-Венан обозначает эту постоянную буквой K . Мы будем называть эту гипотезу "условием пластичности Сен-Венана".

II. Направление наибольшей скорости скольжения совпадает с направлением наибольшего касательного напряжения.

III. Вещество в пластическом состоянии несжимаемо.

Кроме этих уже сформулированных гипотез, Сен-Венан пользуется еще двумя неформулированными гипотезами:

IV. τ_{xz} лежит между главными напряжениями (IX.1.7).

V. Производные по координатам от смещений и скоростей смещений малы.

Перечисленные гипотезы приводят к следующим уравнениям пластического состояния в плоском случае:

$$\tau_{xx/x} + \tau_{xy/y} + \lambda = 0, \quad \tau_{xy/x} + \tau_{yy/y} + \mu = 0; \quad (\text{IX.1.8})$$

$$(\tau_{xx} - \tau_{yy})^2 + 4\tau_{xy}^2 = 4K^2; \quad (\text{IX.1.9})$$

$$(U_{yy} - U_{xx}) : (U_{y/x} + U_{xy}) = (\tau_{yy} - \tau_{xx}) : 2\tau_{xy}; \quad (\text{IX.1.10})$$

$$U_{x/x} + U_{y/y} = 0. \quad (\text{IX.1.11})$$

Уравнения (IX.I.8) суть общие уравнения равновесия механики сплошной среды, уравнение (IX.I.9) есть условие пластичности Сен-Венана, уравнения (IX.I.10) и (IX.I.11) соответствуют гипотезам IV и V.

3^o. Уравнения теории пластичности в трехмерном случае получили Леви [1] на основании условия пластичности Сен-Венана. Леви сохраняет гипотезы I, III и V Сен-Венана. Гипотеза IV становится ненужной, гипотеза II заменяет я следующей:

IV. Дивергент напряжений и тензор скоростей деформаций пропорциональны,

Гипотезы Леви приводят к следующим уравнениям упругого состояния:

$$\begin{aligned} (U_{xx} - U_{yy}) : (\tau_{xx} - \tau_{yy}) &= (U_{yy} - U_{zz}) : (\tau_{yy} - \tau_{zz}) = \\ &= (U_{xy} + U_{yx}) : 2\tau_{xy} = (U_{yz} + U_{zy}) : 2\tau_{yz} = \\ &= (U_{xz} + U_{zx}) : 2\tau_{xz}; \end{aligned} \quad (IX.I.12)$$

$$U_{xxx} + U_{yyy} + U_{zzz} = 0. \quad (IX.I.13)$$

Условие пластичности дает уравнение 0

$$\begin{aligned} K^6 - (a^2 + 3b)K^4/2 + (a^2 + 2b)^2 K^2/16 - (4a^3c + a^2b^2) \\ + 18abc + 4b^3 - 27c^2/64 = 0. \end{aligned} \quad (IX.I.14)$$

Здесь a, b, c - коэффициенты уравнения

$$I + aI^2 - bI - c = \begin{vmatrix} \tau_{xx} - \tau & \tau_{xy} & \tau_{xz} \\ \tau_{xy} & \tau_{yy} - \tau & \tau_{yz} \\ \tau_{xz} & \tau_{yz} & \tau_{zz} - \tau \end{vmatrix} = 0,$$

корни которого суть главные нормали к напряжению.

В уравнениях (IX.I.12-14) следует добавить общие уравнения равновесия (IX.I.3).

Заметим, что в плоском случае из уравнений Леви вытекает гипотеза Сен-Венана гипотеза IV. В самом деле, в плоском

случае $U_x = 0$, $\tau_{xz} = \tau_{yz} = 0$, уравнение (IX.I.13) дает $U_x = -U_y/y$ и из первого уравнения (IX.I.12) следует тогда $\tau_{xy} = (\tau_{xx} + \tau_{yy})/2$. Гипотеза IV вытекает из (IX.I.7) и их последнего неравенства.

Замечание. Уравнение (IX.I.14) дано в книге автора "Основные уравнения математической теории пластичности", изд. АН СССР 1934. В статье Леви [1] дано другое условие, по всей видимости ошибочное, так как из него вытекает условие Сен-Венана для плоской задачи. В сборнике "Теория пластичности" [1] дан перевод другой статьи Леви, в которой дана еще одна запись условия пластичности. Автору кажется, что эта запись также ошибочна, потому что размерностилагаемых в этом условии неодинаковы.

4°. Коротко изложим теорию Мизеса [1]. Существенным ее отличием является иная формулировка условия пластичности. Пусть τ_1 , τ_2 , τ_3 — главные нормальные напряжения. Положим

$$\zeta_1 = \frac{\tau_2 - \tau_3}{2}, \quad \zeta_2 = \frac{\tau_3 - \tau_1}{2}, \quad \zeta_3 = \frac{\tau_1 - \tau_2}{2}; \quad (\text{IX.I.15})$$

так что

$$\zeta_1 + \zeta_2 + \zeta_3 = 0. \quad (\text{IX.I.16})$$

Мизес принимает следующую гипотезу, более общую, чем гипотеза Сен-Венана: в пластическом состоянии напряжения остаются на деле упругости, в пространстве координат ζ_1 , ζ_2 , ζ_3 предельная поверхность упругости представляется в виде замкнутой кривой, лежащей в плоскости (IX.I.16) и обходящей начало координат.

По Сен-Венану $|\zeta_i| \leq K$, $i = 1, 2, 3$. Эти неравенства определяют в пространстве $(\zeta_1, \zeta_2, \zeta_3)$ куб, который пересекается с плоскостью (IX.I.16) по правильному шестиугольнику. Мизес заменил куб Сен-Венана сферой

$$\zeta_1^2 + \zeta_2^2 + \zeta_3^2 = 2K'^2, \quad K' = \text{const}. \quad (\text{IX.I.17})$$

Уравнение (IX.I.17) легко преобразуется к виду

$$\begin{aligned} \tau_i^2 := \frac{1}{6} [(\tau_{xx} - \tau_{yy})^2 + (\tau_{yy} - \tau_{zz})^2 + (\tau_{zz} - \tau_{xx})^2 + \\ + 6(\tau_{xy}^2 + \tau_{yz}^2 + \tau_{xz}^2)] = \frac{4}{3} K'^2. \end{aligned} \quad (\text{IX.I.18})$$

величина τ_i называется интенсивностью касательных напряжений.

В условиях плоской задачи $\tau_i^2 = \frac{1}{4}(\tau_{xx} - \tau_{yy})^2 + \tau_{xy}^2$;
условие пластичности Сен-Венана будет получаться как частный случай условия Мизеса, если положить

$$\mathcal{K}' = \frac{\sqrt{3}}{2} \mathcal{K} ; \quad (\text{IX.I.19})$$

при этом условии Мизеса немного упрощается:

$$\tau_i = \mathcal{K} . \quad (\text{IX.I.20})$$

Остальные уравнения теории пластического состояния по Мизесу совпадают с соответствующими уравнениями Леви.

5°. Вариационный принцип Хаара и Кармана был опубликован в статье [I] этих знаменитых авторов в 1909 г., до появления статьи [I] Мизеса; естественно, что этот принцип был основан на условии пластичности Сен-Венана. Однако, принцип Хаара - Кармана легко формулируется в терминах условия пластичности Мизеса, и в этом случае названный принцип утверждает следующее.

Пусть среда, заполняющая область Ω трехмерного (или двумерного) пространства, находится в упруго-пластическом состоянии. Для простоты допустим, что объемные силы отсутствуют, так что вектор $\mathcal{K} = 0$ и тензор напряжений \mathbb{T} удовлетворяет однородному уравнению

$$\text{div } \mathbb{T} = 0 . \quad (\text{IX.I.21})$$

Пусть $W(\mathbb{T})$ - плотность энергии упругой деформации, выраженная в терминах тензора напряжений \mathbb{T} . Тогда в упруго-пластическом состоянии этот тензор реализует минимум интеграла

$$\mathcal{J}(\mathbb{T}) = \int_{\Omega} W(\mathbb{T}) d\alpha \quad (\text{IX.I.22})$$

на множестве симметричных тензоров \mathbb{T} , удовлетворяющих следующим условиям: а) интеграл (IX.I.22) конечен; б) \mathbb{T} удовлетворяет уравнению (IX.I.21) и известным краевым условиям Коши, означим, что известны внешние силы, приложенные к границе области Ω ; в) \mathbb{T} удовлетворяет с достаточной ограниченностью

$$\tau_i \leq \mathcal{K} , \quad (\text{IX.I.23})$$

где τ_i - интенсивность касательных напряжений тензора \mathbb{T} .

В статье Хаара и Кармана [I] доказывается некоторая общая

теорема, из которой в частности следует: если T решает сформулированную выше одностороннюю вариационную задачу, то в каждой точке области Ω удовлетворяется либо уравнение Эйлера для функционала (IX.1.22), либо уравнение $T_i = \mathcal{H}$ (т.е., условие пластичности).

Следует, однако, иметь в виду, что в пластической зоне должны выполняться также уравнения, содержащие скорости деформаций - так, в трехмерной задаче должны быть удовлетворены уравнения (IX.1.12-13). Поэтому из осуществления решения вариационной задачи Хаара - Кармана еще не вытекает автоматически существование решения задачи об упруго-пластическом состоянии.

Принцип Хаара - Кармана тесно связан с известным принципом Кастильяно о теории упругости. Напомним этот последний принцип: если тензор напряжений реализует минимум функционала (IX.1.22) при условиях а) и б), то этот тензор решает задачу теории упругости при краевых условиях, упомянутых в б). Таким образом, формально можно получить принцип Хаара - Кармана из принципа Кастильяно, если потребовать, чтобы функционал Кастильяно (IX.1.22) достигал минимума на множестве тензоров напряжений, удовлетворяющих не только условиям а), б), но и условию в), т.е., неравенству (IX.1.23).

§ 2. Кручение стержня.

Γ^0 . Пусть нормальное сечение G скручиваемого стержня есть конечная и, в общем случае, многосвязная область плоскости (X, Y) . Если стержень весь находится в упругом состоянии, то отличные от тождественного нуля напряжения $\tau_x := \tau_{xx}$ и $\tau_y := \tau_{yy}$ удовлетворяют уравнениям

$$\tau_{x/x} + \tau_{y/y} = 0, \quad (x, y) \in G; \quad (IX.2.1)$$

$$\tau_x \cos(\nu, x) + \tau_y \cos(\nu, y) = 0, \quad (x, y) \in \partial G; \quad (IX.2.2)$$

$$\iint_G (x\tau_y - y\tau_x) dx dy = m; \quad (IX.2.3)$$

$$\Delta \tau_x = \Delta \tau_y = 0. \quad (IX.2.4)$$

Здесь V - нормаль к ∂G , M - момент внешних сил, приложенных к основанию стержня. Если стержень находится в упруго-пластическом состоянии, то уравнения (IX.2.1-3) сохраняют силу; уравнения (IX.2.4) остаются оправданными в упругой зоне. Приняв, что напряжения в пластической зоне удовлетворяют условию Сен-Венана (уравнение (IX.1.14)) или уравнению Лизеса (IX.1.20), мы должны заменить в пластической зоне уравнения (IX.2.4) одним уравнением

$$\tau_x^2 + \tau_y^2 = \mathcal{K}^2. \quad (IX.2.5)$$

Введем предгильбертово пространство элементы которого суть векторы $\tau = (\tau_x, \tau_y)$, удовлетворяющие уравнению (IX.2.1) и краевому условию (IX.2.2); норму в этом пространстве зададим формулой

$$\|\tau\|^2 = \iint_G (\tau_x^2 + \tau_y^2) dx dy. \quad (IX.2.6)$$

Дополним наше пространство в метрике (IX.2.6), получим гильбертово пространство, которое обозначим через \tilde{H} . Об элементах этого пространства будем говорить, что они удовлетворяют (в обобщенном смысле) уравнениям (IX.2.1-2).

Если стержень находится в упругом состоянии, то в силу принципа Кастильяно задача кручения равносильна задаче о минимуме функционала (IX.2.6) на множестве векторов, удовлетворяющих уравнениям (IX.2.1-3) - иначе говоря, на множестве элементов \tilde{H} , удовлетворяющих уравнению (IX.2.3). Пусть $-q = -(q_x, q_y)$ - какой-нибудь вектор из \tilde{H} , удовлетворяющий этому уравнению. Положим $q + \tau = \sigma = (\sigma_x, \sigma_y)$. Тогда σ удовлетворяет однородному уравнению

$$\iint_G (x\sigma_y - y\sigma_x) dx dy = 0. \quad (IX.2.7)$$

Обозначим через H подпространство тех элементов из \tilde{H} , которые удовлетворяют последнему уравнению. В силу принципа Хара - Кармана решение задачи упруго-пластического кручения является решением следующей односторонней вариационной задачи:

$$\begin{aligned} \|\sigma - q\| &= \min, \sigma \in H; \\ H &= \{\sigma \in \tilde{H}, \|\sigma - q\| \leq \mathcal{K}\}, \end{aligned} \quad (IX.2.8)$$

$$\|\tau\| = \sup_{(x,y) \in G} \operatorname{ess} \sqrt{\tau_x^2 + \tau_y^2}.$$

2°. Исследуем множество M , определенное в формулах (IX.2.8). При больших по модулю значениях M оно будет пустым. Действительно,

$$|M| \leq \iint_G \rho \sqrt{\tau_x^2 + \tau_y^2} dx dy; \quad \rho^2 = x^2 + y^2;$$

в упруго-пластическом состоянии $\tau_x^2 + \tau_y^2 \leq K^2$, поэтому

$$|M| \leq K \iint_G \rho dx dy. \quad (\text{IX.2.9})$$

Центральная часть последнего неравенства есть некоторая постоянная: если в уравнении (IX.2.8) $|M|$ больше этой постоянной, то M пусто.

Назовем момент M допустимым, если ему соответствует непустое множество M . Изменив в случае необходимости ориентацию осей координат, можно считать, что $\partial \bar{M} > 0$. Обозначим точную верхнюю границу допустимых моментов через \bar{M} и докажем, что множество допустимых моментов есть либо отрезок $0 \leq M \leq \bar{M}$, либо полуинтервал $0 \leq M < \bar{M}$. Достаточно доказать, что последний полуинтервал содержится в множестве допустимых моментов.

В любой лезвой полускрестности точки \bar{M} содержится хотя бы один допустимый элемент M' . Ему соответствует по крайней мере один вектор $\tau' = (\tau'_x, \tau'_y)$, такой, что

$$\iint_G (x\tau'_y - y\tau'_x) dx dy = M'; \quad \|\tau'\| \leq K.$$

Пусть $0 \leq M'' < M'$. Поиском $\tau'' = M'' \tau' / M'$. Тогда $\tau'' \in \bar{M}$ и

$$\iint_G (x\tau''_y - y\tau''_x) dx dy = M'', \quad \|\tau''\| = \frac{M''}{M'} \|\tau'\| < K.$$

Последнее соотношение означает, что множество \bar{M} , которое соответствует моменту M'' , содержит элемент τ'' и потому не пусто. Этим и доказано, что полуинтервал $0 \leq M < \bar{M}$ содержится в множестве допустимых моментов.

Докажем теперь, что множество M , определенное формулами (IX.2.8), замкнуто в норме $|\cdot|$. Пусть $\tau_n = (\tau_{nx}, \tau_{ny})$ - последовательность векторов из M , и пусть $\|\tau_n - \tau_0\| \xrightarrow{n \rightarrow \infty} 0$.

Докажем, что $\tau_0 \in M$. Из сходящейся в H последовательности $\{\tau_n\}$ можно выделить подпоследовательность $\{\tau_{n_k}\}$, сходящуюся к τ_0 почти всюду в G . Так как $\|\tau_{n_k}\| = \sup \sqrt{\tau_{n_k x}^2 + \tau_{n_k y}^2} < K$, то и $\sup \sqrt{\tau_{\alpha x}^2 + \tau_{\alpha y}^2} < K$ почти всюду в Ω . Это и означает, что $\tau_0 \in M$ и, следовательно, множество M замкнуто.

Из известных общих теорем (см. § I гл. VIII) следует, что задача (IX.2.8) имеет одно и только одно решение.

3°. Задача (IX.2.8) есть частный случай задачи (VIII.1.1-2), в которой, в частности, следует положить $d = K$. Если $0 \leq m < \bar{m}$, то, как мы видели в п.2°, в множестве M существует элемент, полунорма которого меньше K . Примем этот элемент за q . На основании результатов § 5 гл. VIII можно утверждать следующее.

Решение задачи (IX.2.8) устойчиво относительно малых возмущений вектора q и постоянной пластичности K , если крутящий момент лежит в полуинтервале $0 \leq m < \bar{m}$.

Пусть момент m заменен моментом m' , мало отличающимся от m и по-прежнему лежащим в полуинтервале $[0, \bar{m})$. Если моменту m соответствует вектор q , $\|q\| < K$, то моменту m' можно привести в соответствие вектор $q' = m'q/m$; если m и m' достаточно близки, то $\|q'\| < K$ и величина $\|q - q'\|$ достаточно мала. Из сказанного вытекает следующее утверждение:

Решение задачи (IX.2.8) устойчиво относительно малых возмущений крутящего момента, не выходящих из полуинтервала $[0, \bar{m})$.

Ниже мы покажем, что при том же условии $0 \leq m < \bar{m}$ решение задачи (IX.2.8) устойчиво относительно малых изменений области G .

4°. Для последующих ссылок приведем некоторые известные факты теории упругого кручения стержней.

Пусть область G - $(n+1)$ -связная. Тогда ∂G состоит из $n+1$ связных компонент, которые мы обозначим через ∂G_j , $j=0, 1, \dots, n$; через ∂G_0 обозначим ту компоненту, которая ограничивает область G извне. Обозначим еще через G_j , $j=1, 2, \dots, n$, ту область, которая извне ограничена контуром ∂G_j . Пусть $G_0 = G \cup G_1 \cup \dots \cup G_n$; G_0 есть конечная односвязная область, ограниченная контуром ∂G_0 . Можно определить G_0 как "меньшую" односвязную область, содержащую область G .

Уравнение (IX.2.1) позволяет ввести так называемую функцию напряжений Прандтля

$$\tau_x = \frac{\partial u}{\partial y}, \quad \tau_y = -\frac{\partial u}{\partial x}. \quad (\text{IX.2.1})$$

Из краевого условия (IX.2.2) следует, что на границе области функция Прандтля удовлетворяет условию

$$u|_{\partial G_j} = A_j = \text{const}. \quad (\text{IX.2.1})$$

При изменении функции Прандтля на постоянное слагаемое напряжения не меняются, поэтому можно положить $A_0 = 0$.

Как видно из уравнений (IX.2.4), в упругой области функция Прандтля удовлетворяет уравнению Пуассона

$$-\Delta u = C_0, \quad C_0 = \text{const}. \quad (\text{IX.2.1})$$

Известна формула, связывающая крутящий момент M с функцией Прандтля $u(x, y)$, именно

$$M = 2 \iint_G u dx dy + 2 \sum_{j=1}^n A_j |G_j|. \quad (\text{IX.2.1})$$

Продолжим функцию u на всю область G_0 , положив $u(x, y) = 0$ ($(x, y) \in G_j$). Тогда формула (IX.2.13) принимает более простой вид

$$M = 2 \iint_{G_0} u dx dy. \quad (\text{IX.2.1})$$

5°. Задачу кручения можно свести к вариационной задаче, которая по форме отлична от задачи (IX.2.8). Если сжимаемый стержень находится в упругом состоянии, то задачу (IX.2.8) можно записать в терминах функции Прандтля следующим образом:

$$\iint_G |\nabla u|^2 dx dy = \min,$$

$$u = \begin{cases} 0 & \text{на } \partial G_0, \\ A_j = \text{const} & \text{на } \partial G_j, \quad 1 \leq j \leq n, \end{cases} \quad (\text{IX.2.1})$$

$$2 \iint_G u dx dy + 2 \sum_{j=1}^n A_j |G_j| = M.$$

Это задача на условный минимум; по теореме Эйлера существует такая постоянная C_0 , что

$$\delta \left\{ \iint_G [|\nabla w|^2 - 2C_0 u] dx dy - 2C_0 \sum_{j=1}^n A_j |G_j| \right\} = 0, \quad (\text{IX.2.16})$$

δ - знак вариации; функция u должна удовлетворять краевому условию $\#z$ (IX.2.15). Если эту функцию продолжить на G_0 , как об этом сказано в п.4^o, то уравнение (IX.2.16) примет вид

$$\delta \iint_{G_0} [|\nabla w|^2 - 2C_0 u] dx dy = 0. \quad (\text{IX.2.17})$$

Обозначим через \bar{H}_0 подпространство пространства $\dot{W}_2^{(1)}(G_0)$, образованное функциями, которые остаются постоянными в каждой из внутренних областей G_j . Норму в $\dot{W}_2^{(1)}(G_0)$ зададим формулой

$$\|u\|^2 = \iint_{G_0} |\nabla u|^2 dx dy;$$

соответственно скалярное произведение равно интегралу

$$[u, v] = \iint_{G_0} \nabla u \cdot \nabla v dx dy.$$

Вариационная задача (IX.2.15) приводится к виду

$$|u - u_0| = \min, \quad u \in \bar{H}_1,$$

где u_0 , в соответствии с теоремой Ф.Рисса, определяется тождеством

$$\forall u \in \bar{H}_0, \quad [u, u_0] = C_0 \iint_{G_0} u dx dy.$$

Применяя принцип Хэара - Кармана, мы сведем задачу упруго-пластического кручения к следующей односторонней вариационной задаче:

$$|u - u_0| = \min, \quad u \in \bar{M}; \quad (\text{IX.2.18})$$

$$\bar{M} = \{u \in \bar{H}_0 : \|u\| \leq K\};$$

$$\|u\| = \sup_G |\nabla u|$$

Задача (IX.2.18) сформулирована в статье Аллен Ланшон [1]. При любом конечном значении C_0 эта задача имеет одно и только одно решение, которое мы обозначим через $u_* = u_*(x, y, C_0)$. В статье Ланшон [1] доказывается, что решение задачи (IX.2.18) удовлетворяет всем дифференциальным уравнениям и крайним условиям задачи упругопластического кручения.

Если задан крутящий момент M , то C_0 следует определять из уравнения

$$2 \iint_{G_0} u_*(x, y, C_0) dx dy = M. \quad (\text{IX.2.19})$$

Каждому конечному значению C_0 соответствует по формуле (IX.2.19) одно и только одно допустимое значение M . Обратным образом, как видно из результатов пп. 2^о - 4^о, каждому допустимому значению M из полуинтервала $0 < M < \bar{M}$ соответствует одна и только одна функция Прандтля, которая в упругой области удовлетворяет уравнению вида (IX.2.12), иначе говоря, каждому допустимому значению M из полуинтервала $[0, \bar{M})$ соответствует одно и только одно конечное значение C_0 , и при любом значении $M \in [0, \bar{M})$ уравнение (IX.2.19) имеет одно и только одно решение.

6^о. Займемся вопросом об устойчивости решения задачи (IX.2.18) относительно малых возмущений области G . Наряду с задачей (IX.2.18) рассмотрим такую же задачу для отержня из того же материала и с нормальным сечением G' , близким к G . Близость областей G и G' определим следующим образом. Обозначим через ξ, η и X, Y декартовы координаты в областях G' и G соответственно и допустим, что существует преобразование координат S , которое взаимно однозначно переводит G' в G и G_0 в G (здесь G_0 - наименьшая односвязная область, содержащая область G'), причем якобиана матрица преобразования S имеет вид $I + O(\varepsilon)$, где I - единичная матрица, а ε - малое положительное число. Тогда, если T - якобиево преобразование S , то $T = I + O(\varepsilon)$.

Задачу кручения для области G' можно записать так:

$$\iint_{G_0} \left[\left(\frac{\partial u'}{\partial \xi} \right)^2 + \left(\frac{\partial u'}{\partial \eta} \right)^2 - 2C_0 u' \right] d\xi d\eta = \min; \quad u' \in \bar{H}_0,$$

$$\min_{G_0} \int_{G_0} \left[\left(\frac{\partial u'}{\partial \xi} \right)^2 + \left(\frac{\partial u'}{\partial \eta} \right)^2 \right] \leq M^2; \quad (\text{IX.2.20})$$

здесь \bar{H}_0 - подпространство пространства $W_2^{(1)}(G_0)$, образованное

функциями, поставленным в каждой из связанных компонент множества $G_0 \setminus G'$.

В (IX.2.20) выполним преобразование S ; заметим, что оно переводит \bar{H}'_0 в \bar{H}_0 . Обозначая $\nabla' = (\frac{\partial}{\partial \xi}, \frac{\partial}{\partial \eta})$, получаем

$$\iint_{G_0} \{ (u_{/x})^2 |\nabla' x|^2 + 2u_{/x} u_{/y} (\nabla' x, \nabla' y) + (u_{/y})^2 |\nabla' y|^2 - 2c_0 u \} J dx dy = \min;$$

$$u(x, y) = u(\xi, \eta); \quad u \in \bar{H}_0;$$

$$\sup_{G_0} \{ (u_{/x})^2 |\nabla' x|^2 + 2u_{/x} u_{/y} (\nabla' x, \nabla' y) + (u_{/y})^2 |\nabla' y|^2 \} \leq \mathcal{K}^2. \quad (\text{IX.2.21})$$

В пространстве \bar{H}_0 введем новую норму $|\cdot|_1$ (и соответственно новое скалярное произведение $[\cdot, \cdot]$), положив

$$|u|_1^2 = \iint_{G_0} \{ (u_{/x})^2 |\nabla' x|^2 + 2u_{/x} u_{/y} (\nabla' x, \nabla' y) + (u_{/y})^2 |\nabla' y|^2 \} J dx dy. \quad (\text{IX.2.22})$$

Функционал $l_1 u = \iint_{G_0} c_0 J u dx dy$ ограничен в норме (IX.2.22), поэтому $l_1 u = [u, u_1]_1$, где u_1 — некоторый элемент из \bar{H}_1 . Наконец, обозначим через \bar{M} множество функций из \bar{H}_0 , удовлетворяющих последнему из соотношений (IX.2.21) и через $\|u\|_1^2$ — левую часть этого соотношения. Задача (IX.2.21) принимает вид

$$|u - u_1|_1 = \min, \quad u \in \bar{M};$$

$$\bar{M} = \{ u \in \bar{H}_0 : \|u\|_1 \leq \mathcal{K} \}. \quad (\text{IX.2.2})$$

7°. При переходе от задачи (IX.2.18) к задаче (IX.2.23) меняются все параметры, определяющие задачу: элемент c_0 , норма $|\cdot|$ и множество \bar{M} . Покажем, что все эти изменения малы в смысле определения § 5 гл. VIII.

Прежде всего докажем, что нормы $|\cdot|$ и $|\cdot|_1$ связаны неравенством вида (VIII.5.3). Выражение под знаком интеграла в (IX.2.22) есть квадратичная форма относительно $u_{/x}$ и $u_{/y}$; нетрудно видеть, что ее собственные числа имеют вид $\lambda_k = 1 + O(\varepsilon)$, $k=1, 2$. В п. 6° было отмечено, что $\gamma = 1 + O(\varepsilon)$. Из сказанного следует существование такой постоянной δ_1 , зависящей только от ε , что $\forall u, \delta_1 = 0$ и $(1 - \delta_1)^2 \leq \gamma \lambda_k \leq (1 + \delta_1)^2$; $k=1, 2$. Тогда $\varepsilon \rightarrow 0$

$$(1 - \delta_1)|u| \leq |u_1| \leq (1 + \delta_1)|u|.$$

Теперь оценим величину $|u_0 - u_1|$. Рассмотрим разность $[u, u_1] - [u, u_1]_1 =$

$$= \iint_{G_0} \left\{ u_{1/x} u_{1/x} (1 - \gamma |\nabla' x|^2) + u_{1/y} u_{1/y} (1 - \gamma |\nabla' y|^2) - (u_{1/x} u_{1/y} + u_{1/y} u_{1/x}) (\gamma x \cdot \gamma y) \right\} dx dy$$

Из последнего тождества находим $[u, u_1] - [u, u_1]_1 = |u| \cdot |u_1| \cdot O(\varepsilon)$ и, следовательно,

$$[u, u_0 - u_1] = [u, u_0] - [u, u_1] + |u| \cdot |u_1| \cdot O(\varepsilon).$$

Далее, принимая, что $c_0 - c'_0 = O(\varepsilon)$, находим

$$\begin{aligned} [u, u_0] - [u, u_1] &= c_0 \iint_{G_0} u dx dy - c'_0 \iint_{G_0} u dx dy = \\ &= (c_0 - c'_0) \iint_{G_0} u dx dy + c'_0 \iint_{G_0} u (1 - \gamma) dx dy = |u| \cdot O(\varepsilon) \end{aligned}$$

и $[u, u_0 - u_1] = |u| (1 + |u_1|) O(\varepsilon)$. Полагая здесь $u = u_0 - u_1$, находим

$$|u_0 - u_1| \leq C \varepsilon (1 + |u_1|), \quad C = \text{const}. \quad (\text{IX.2.24})$$

Отсюда $|u_1| - |u_0| \leq C \varepsilon (1 + |u_1|)$ и $|u_1| \leq |u_0| + O(\varepsilon)$. Теперь то же неравенство (IX.2.24) дает $|u_0 - u_1| = O(\varepsilon)$; существует, следовательно, такая величина δ_2 , зависящая только от ε , что $\lim_{\varepsilon \rightarrow 0} \delta_2 = 0$ и $|u_0 - u_1| \leq \delta_2$.

8°. Исследуя влияние замены множества \bar{M} на \bar{M}' , примем пока, что $G' \subset G$. Введенные выше пространства \bar{H}_0 и \bar{H}'_0 будем теперь обозначать через $H_0(G)$ и $H_0(G')$. Функции класса $\bar{H}_0(G')$ определим в $G_0 \setminus G'_0$, положив их там равными нулю; после этого $H_0(G')$ станет подпространством пространства $H_0(G)$.

Выше мы определили S как преобразование координат; расширим это определение, перенеся его также на функции, заданные в соответствующих областях, так что если функция $u(\xi, \eta)$ определена в G_0 , то $(Su)(x, y) = u(\xi(x, y), \eta(x, y))$, где $(x, y) = S(\xi, \eta)$. При таком определении S переводит \bar{M}' в \bar{M} ; докажем, что при этом выполняется неравенство (VII.5.9). Пусть $u \in \bar{M}$. При этом

$$S^{-1}u = \begin{cases} u(\xi(x,y), \eta(x,y)), & \xi(x,y) \in G_0, \\ 0 & \xi(x,y) \in G_0 \setminus G_0' \end{cases}$$

и

$$|u - S^{-1}u| = \iint_{G_0 \setminus G_0'} |\nabla u|^2 dx dy + \iint_{G_0'} [\nabla(u(\xi, \eta) - (S^{-1}u)(\xi, \eta))]^2 d\xi d\eta.$$

Первый интеграл справа, в силу условия пластичности Сен-Венана, имеет порядок $O(\varepsilon)$. Второй интеграл равен

$$\iint_{G_0'} [(u_{\xi}^2 [(1-\xi_x)^2 + \xi_y^2] + 2u_{\xi} u_{\eta} [(1-\xi_x)\eta_y + \xi_y(1-\eta_x)] + (u_{\eta}^2 [\eta_x^2 + (1-\eta_y)^2])] dx dy.$$

Якобиева матрица преобразования S имеет вид $I + O(\varepsilon)$; отсюда следует, что последний интеграл оценивается величиной $|u|^2 O(\varepsilon^2)$. Окончательно, $|u - S^{-1}u|^2 = O(\varepsilon) + |u|^2 O(\varepsilon^2)$, $|u - S^{-1}u| = O(\sqrt{\varepsilon}) + |u| O(\varepsilon) = (1 + |u|) O(\sqrt{\varepsilon})$. Это значит, что существует величина δ_3 , зависящая только от ε , такая, что $\lim_{\varepsilon \rightarrow 0} \delta_3 = 0$ и $|u - S^{-1}u| < \delta_3(1 + |u|)$.

Из доказанного в настоящем пункте следует, что задача упруго-пластического кручения устойчива относительно малых возмущений области сечения стержня, если возмущенная область лежит внутри первоначальной.

9°. Пусть область G' получена малым возмущением области G , но G' не вкладывается в G . Построим область \tilde{G} , которая может быть получена малым возмущением каждой из областей G и G' , и притом такую, что как G , так и G' вкладывается в \tilde{G} . По доказанному в п. 8°, функции Прандля u и u' , соответствующие сечениям G и G' , мало отличаются от \tilde{u} - функции Прандля для сечения \tilde{G} . Но тогда u и u' мало различаются между собой.

Таким образом, решение задачи упруго-пластического кручения устойчиво относительно малых возмущений сечения стержня, если малость возмущения понимается в смысле п. 6° настоящего параграфа.

10°. На решение задачи упруго-пластического кручения естественным образом распространяются результаты п. 2 и § 6 гл. V. Таким образом, если приближенно значение \tilde{u} функции Прандля

$u_*(x, y)$ определено по методу Рунца (в частности, по методу конечных элементов) как элемент соответствующего подпространства H_n , а $u_*(x, y)$ есть точные значения функции Прандтля, то

$$\|u_* - u_r^{(n)}\| = O(\sqrt{\epsilon_n(u_*)}); \quad (IX.2.25)$$

$$C_n(u) = \inf_{u^{(n)} \in H_n} [\|u - u^{(n)}\| + \| \|u - u^{(n)}\| \|],$$

причем норма $\|\cdot\|$ и полунорма $\|\cdot\|$ соответственно определяются формулами

$$\|u\|^2 = \iint_G |\nabla u|^2 dx dy, \quad \| \|u\| \| = \sup_G |\nabla u|.$$

Погрешность искажения приближенного решения $u_*(x, y)$ имеет вид $O(\sqrt{\|\Gamma_n\|} + \|\delta^{(n)}\|)$, где Γ_n - искажение матрицы Рунца, $\delta^{(n)}$ есть искажение столбца свободных членов системы Рунца.

II°. В литературе много внимания уделяется приближенному решению задачи (IX.2.18). В случае односвязной области этот вопрос обстоятельно изложен в книге Гловинского, Лисеня и Тремольера [1]: некоторые результаты содержатся в более ранней книге Черноусько и Бенничука [1]. Случаю много связного сечения посвящена работа Гловинского и Ланшон [1]. Оценки погрешности приближения в этих работах отсутствуют. Эластическое кручение отержней с односвязным сечением изучалось как в точной, так и в приближенной постановке в ряде работ и независимо от вариационного принципа. Первой в этом ряду была знаменитая работа Надаи (см. Надаи [1]). Интересные результаты были получены в этом направлении Галиным [1 - 2]. Отметим еще работу Перлина [1].

§ 3. Плоская задача

I°. Введем следующие обозначения: σ - постоянная Пуассона упругой среды,

$$T = \begin{pmatrix} \tau_{xx} & \tau_{xy} \\ \tau_{xy} & \tau_{yy} \end{pmatrix} \quad (IX.3.1)$$

- тензор напряжений,

$$\tau_m = \left[\frac{1}{4} (\tau_{xx} - \tau_{yy})^2 + \tau_{xy}^2 \right]^{1/2} \quad (IX.3.2)$$

- максимальное касательное напряжение, X и Y - составляющие объемных сил, X_y и Y_x - составляющие поверхностных усилий. Далее, G - конечная область в плоскости (x, y) , $\partial G = \bigcup_{j=0}^n \partial G_j$

граница области G , ∂G_j - ее связные компоненты. Тензор напряжений удовлетворяет внутри G уравнениям равновесия

$$\tau_{xx,x} + \tau_{xy,y} + X = 0, \quad \tau_{xy,x} + \tau_{yy,y} + Y = 0, \quad (IX.3.3)$$

а на ∂G - крайним условиям

$$\tau_{xx} \cos(\nu, x) + \tau_{xy} \cos(\nu, y) = X_\nu, \quad \tau_{xy} \cos(\nu, x) + \tau_{yy} \cos(\nu, y) = Y_\nu. \quad (IX.3.4)$$

Если среда находится в упругом состоянии, то по принципу Кастильяно тензор упругих напряжений, который мы обозначим через \hat{Q} , сообщает минимальное значение интегралу

$$\iint_G [(1-2\sigma)(\tau_{xx} + \tau_{yy})^2 + 4\tau_m^2] dx dy \quad (IX.3.5)$$

на множестве тензоров, удовлетворяющих уравнениям (IX.3.3-4). По условию пластичности Сен-Венана, в пластическом состоянии в условиях плоской задачи выполняется равенство

$$\tau_m^2 = \frac{1}{4} [(\tau_{xx} - \tau_{yy})^2 + 4\tau_{xy}^2] = \mathcal{K}^2, \quad (IX.3.6)$$

а тогда, по принципу Хаара - Кармана, тензор T_x упруго-пластических напряжений решает одностороннюю задачу о минимуме интеграла (IX.3.5) на множестве тензоров, удовлетворяющих уравнениям (IX.3.3-4) и, кроме того, неравенству

$$(\tau_{xx} - \tau_{yy})^2 + 4\tau_{xy}^2 \leq 4\mathcal{K}^2. \quad (IX.3.7)$$

2°. Сформулированную здесь одностороннюю вариационную задачу сведем к задаче вида (УП. I. i-2). Введем гильбертово пространство \hat{H} тензоров вида (IX.3.1), составляющие которых квадратично суммируемы в L_2 , с нормой

$$|\hat{T}| = \left\{ \iint_G [(1-2\sigma)(\tau_{xx} + \tau_{yy})^2 + 4\tau_m^2] dx dy \right\}^{1/2} \quad (IX.3.8)$$

и с соответствующим скалярным произведением, а также подпространство \hat{H} этого пространства, образованное тензорами, удовлетворяющими однородным уравнениям равновесия

$$\tau_{xx,x} + \tau_{xy,y} = 0, \quad \tau_{xy,x} + \tau_{yy,y} = 0 \quad (IX.3.9)$$

и однородным крайним условиям

$$\tau_{xx} \cos(\nu, x) + \tau_{xy} \cos(\nu, y) = 0, \quad \tau_{xy} \cos(\nu, x) + \tau_{yy} \cos(\nu, y) = 0. \quad (IX.3.10)$$

Аналогично тому, как это сделано в книге автора [5] для трехмерной задачи теории упругости, можно доказать, что пространство $\mathcal{H} \in \mathcal{H}$ состоит из тензоров, которые удовлетворяют соотношениям (I.3.3-4) при подходящих X, Y, X_v, Y_v и уравнениям плоской задачи теории упругosti; в частности, $\hat{g} \in \mathcal{H} \in \mathcal{H}$.

Обозначим через $\| \cdot \|$ полуnormу

$$\| T \| = \sup_{(x,y) \in G} \tau_m. \quad (\text{IX.3.II})$$

Положим $\bar{g} = \hat{g} + V$, где V - произвольный тензор из \mathcal{H} . Если $\inf_{V \in \mathcal{H}} \| \bar{g} \| > \mathcal{K}$, то односторонняя вариационная задача $p. 1^{\circ}$ неразрешима (см. гл. VIII, у I, п. 2^o). Будем считать, что $\inf_{V \in \mathcal{H}} \| \bar{g} \| < \mathcal{K}$. Тогда найдется такой тензор $V_0 \in \mathcal{H}$, что $\| V - g \| < \mathcal{K}$, где $g = \hat{g} + V_0$. Обозначая $T = V - g, V \in \mathcal{H}$, мы приведем задачу $p. 1^{\circ}$ к виду

$$\begin{aligned} |V - g| &= \min, V \in \mathcal{H}; \\ \mathcal{M} &= \{V \in \mathcal{H} : \| V - g \| \leq \mathcal{K}\}. \end{aligned} \quad (\text{IX.3.I2})$$

Последнюю задачу можно слегка упростить. Имеем $V - g = (V - V_0) - \hat{g}$; слагаемые справа ортогональны в метрике (IX.3.8) и $|V - g|^2 = |V - V_0|^2 + |\hat{g}|^2$. Постоянная величина $|\hat{g}|^2$ не влияет на минимизирующий элемент V , и ее можно отбросить. В результате задача (IX.3.I2) заменяется следующей:

$$\begin{aligned} |V - V_0| &= \min, V \in \mathcal{H}; \\ \mathcal{M} &= \{V \in \mathcal{H} : \| V - g \| \leq \mathcal{K}\}. \end{aligned} \quad (\text{IX.3.I3})$$

Вместо того, как это сделано в п. 2^o §2, можно доказать, что множество \mathcal{M} задачи (IX.3.I2) замкнуто в норме $\| \cdot \|$.

Если граница области достаточно гладкая, то основные задачи линейной теории можно свести к интегральным уравнениям, ф. д. вдоль ее или контурными; подробнее об этом см. книги Мухомелишвили [1] автора [1,6], Купрадзе и др. [1]. Несложный анализ этих уравнений показывает, что при подходящем выборе необходимых матриц тензоров напряжений, регулярный вторую краевую задачу теории упругости, устойчив в метрике \mathcal{C} по отношению к малым изменениям объемных и поверхностных сил, а также области, занятой упругой средой.

Опираясь на этот факт, мы сделаем следующее допущение: при

отсюда теперь легко заключить, что такой же uniqueness обладает тензоры g и V_0 . Из результатов § 5 гл. VIII следует, что тензор напряжений, решающий плоскую задачу теории упруго-пластического состояния, устойчив относительно малых возмущений области, если возмущенная область $G' \subset G$. Последнее требование несущественно - от него можно освободиться так же, как в п. 9° § 2.

4°. Для плоской задачи упруго-пластического состояния верны оценки погрешностей аппроксимации и искажения, аналогичные оценкам п. 10° § 2. Мы не будем здесь останавливаться на этом подробнее.

5°. Поставим вопрос: при каких обстоятельствах можно утверждать, что решение плоской задачи упруго-пластического состояния, полученное из принципа Хаара - Кармана, является на самом деле решением этой задачи. По самой постановке задачи, тензор T упругих напряжений, полученный на основе этого принципа, удовлетворяет уравнениям равновесия (IX.3.9) (мы принимаем для простоты, что объемные силы отсутствуют). Далее, как показано в работе Хаара и Кармана [1], в упругой зоне этот тензор удовлетворяет уравнению Эйлера, которое в данном случае сведется к дифференциальному уравнению

$$\Delta(\tau_{xx} + \tau_{yy}) = 0 \quad (\text{IX.3.16})$$

и к краевому условию (если границы области Ω и упругой зоны имеют общую часть)

$$T \cdot \nu = k_\nu; \quad (\text{IX.3.17})$$

здесь ν - внешняя нормаль к ∂G , k_ν - вектор напряжений, действующих на "элементарную площадку" границы ∂G . В силу результатов той же статьи [1] Хаара и Кармана, в пластической зоне тензор T удовлетворяет условию пластичности Сен-Венана (IX.I.9).

Таким образом, в упругой зоне три оставшихся тензора T удовлетворяют системе трех дифференциальных уравнений (IX.3.9) и (IX.3.16); как известно, если T удовлетворяет этой системе, то можно найти такой вектор смещений, что соответствующий тензор деформаций вместе с тензором T удовлетворяет уравнениям закона Гука.

В пластической зоне тензор T удовлетворяет системе трех

уравнений (IX.3.9) и (IX.1.9). Этот тензор будет решать плоскую задачу упруго-пластического состояния, если в пластической зоне существует вектор смещений, который удовлетворяет уравнениям (IX.1.10-II) и вместе со своей производной по времени непрерывно переходит в вектор упругих смещений при переходе через общую границу упругой и пластической зон (см. Собслев [1]). Можно указать простое достаточное условие, при котором последние требования выполнены, так что принцип Хаара-Кармана приводит к решению плоской задачи упруго-пластического состояния. Это будет, если данные задачи не зависят от времени. В этом случае тензор напряжений, а также вектор смещений в упругой зоне не зависят от времени: на границе раздела зон упругости и пластичности вектор скоростей смещений равен нулю. Но тогда можно удовлетворить уравнениям (IX.1.10-II), положив вектор скоростей смещений \dot{U} тождественно равным нулю в пластической зоне.

Что касается вектора смещений в пластической зоне, то он не определяется из уравнений Сен-Венана и его можно подчинить только двум требованиям: он непрерывно переходит в вектор упругих смещений на общей границе обеих зон; он не зависит от времени.

Заметим, что высказанные здесь соображения в полной мере относятся и к трехмерной задаче упруго-пластического состояния (см. ниже, § 4).

6°. В заключение данного параграфа заметим следующее. Полуорма (IX.3.II) не является непрерывно дифференцируемой, и построение приближенного решения наталкивается на некоторые трудности. Можно упростить задачу, заменив полуорму (IX.3.II) следующей

$$\begin{aligned} \|T\|_p &= \frac{1}{|G|^{1/2p}} \|t_m\|_{q,p} = \\ &= \frac{1}{|G|^{1/2p}} \left\{ \iint_G \left[\frac{1}{4} (\tau_{xx} - \tau_{yy})^2 + 4\tau_{xy}^2 \right]^p d\alpha d\beta \right\}^{1/2p} \end{aligned} \quad (IX.3.I8)$$

а неравенство $\|V - \bar{q}\| \leq \mathcal{K}$ - неравенством

$$\|V - \bar{q}\|_p \leq \mathcal{K}. \quad (IX.3.I9)$$

В формуле (IX.3.I8) p достаточно большое натуральное число; очевидно, новая полуорма непрерывно дифференцируема. Обозначим множество тензоров V , удовлетворяющих неравенству (IX.3.I9), через \mathcal{M}_p . Так как $\| \cdot \|_p \leq \| \cdot \|$, то множество \mathcal{M}_p не пусто, если не пусто множество \mathcal{M} . Далее, очевидно, что множество \mathcal{M}_p выпук-

люе. Докажем, что оно замкнуто в норме $|\cdot|_p$. Пусть $\{V_n\} \subset {}^1V_p$ — последовательность сходящаяся в упомянутой норме (формула (IX.3.8)) к некоторому тензору V_0 . Тогда $f_n \rightarrow f_0$ и $g_n \rightarrow g_0$ в $L_2(G)$, где

$$f_n = V_{n\alpha\alpha} - V_{n\gamma\gamma}, \quad g_n = V_{n\alpha\gamma},$$

$$f_0 = V_{0\alpha\alpha} - V_{0\gamma\gamma}, \quad g_0 = V_{0\alpha\gamma}.$$

Положим еще $g_{\alpha\alpha} - g_{\gamma\gamma} = \tilde{f} - \tilde{g}$, тогда в $L_2(G)$

$$\begin{matrix} f_n - \tilde{f} \\ g_n - \tilde{g} \end{matrix} \rightarrow \begin{matrix} f_0 - \tilde{f} \\ g_0 - \tilde{g} \end{matrix}, \quad g_n - \tilde{g} \rightarrow g_0 - \tilde{g}. \quad (\text{IX.3.20})$$

Отсюда вытекает существование такой подпоследовательности $\{V_{n_k}\}$, что те же соотношения (IX.3.22) для этой подпоследовательности выполняются почти всюду в G . По теореме Вату (см., напр. мер. Вулик [1])

$$\|V_0 - g_0\|_p = \frac{1}{|G|^{1/p}} \left\{ \iint_G \left[\frac{1}{4} (x_0 - \tilde{f})^2 + (y_0 - \tilde{g})^2 \right]^p dx dy \right\}^{1/2p} <$$

$$< \frac{1}{|G|^{1/2p}} \sup_{n_k} \left\{ \iint_G \left[\frac{1}{4} (V_{n_k} - \tilde{f})^2 + (y_{n_k} - \tilde{g})^2 \right]^p dx dy \right\}^{1/2p} =$$

$$= \sup_{n_k} \|V_{n_k} - g\|_p < K.$$

Это означает, что $V_0 \in {}^1V_p$ и, следовательно, множество замкнуто.

Аналогично можно упростить задачи § 2 и § 4 (см. ниже).

§ 4. Трехмерная задача

Γ^n . Пусть известны векторы \tilde{z} — осевых и \tilde{z}_ν — поверхностных сил, действующих на упруго-пластическую среду, которая занимает конечную область Ω трехмерного пространства. Тензор напряжений удовлетворяет в Ω дифференциальному уравнению

$$\operatorname{div} T + \tilde{z} = 0, \quad (\text{IX.4.1})$$

а на границе этой области — краевому условию

$$T \cdot \nu_{|\partial\Omega} = \tilde{z}_\nu; \quad (\text{IX.4.2})$$

ν — внешняя норма к Ω .

Если среда находится в упругом состоянии, то по принципу Кастильяно тензор T сообщает минимальное значение функционалу.

$$\iint_{\Omega} \left[\frac{1-2\nu}{6(1+\nu)} (\tau_{xx} + \tau_{yy} + \tau_{zz})^2 + \tau_i^2 \right] dx dy dz \quad (\text{IX.4.3})$$

на множестве тензоров, удовлетворяющих уравнениям (IX.4.1-2) через τ_i обозначена интенсивность касательных напряжений. Применяя принцип Хаара-Кармана, будем считать, что тензор напряжений в упруго-пластическом состоянии сообщает минимальное значение функционалу (IX.4.3) на множестве тензоров, которые удовлетворяют, кроме уравнений (IX.4.1-2), еще и нечетности (IX.1.23).

2°. Только что сформулированную одностороннюю вариационную задачу преобразуем к обычному виду. Примем, что тензор напряжений квадратично суммируем в Ω . На множестве таких тензоров введем норму

$$\|T\|^2 = \iint_{\Omega} \left[\frac{1-2\nu}{6(1+\nu)} (\tau_{xx} + \tau_{yy} + \tau_{zz})^2 + \tau_i^2 \right] dx dy dz, \quad (\text{IX.4.4})$$

и обозначим полученное таким образом гильбертово пространство через \tilde{H} . Выделим множество H элементов пространства \tilde{H} , которые удовлетворяют однострочным уравнениям

$$\text{div } T = 0, \quad T \cdot \nu|_{\partial\Omega} = 0. \quad (\text{IX.4.5})$$

Как доказано в книге автора [3] H есть подпространство в \tilde{H} . Ортогональное к H подпространство состоит из тензоров упругих напряжений: это значит, что если $T \in H \perp H$, то T удовлетворяет уравнениям (IX.4.1-2) при подходящем выборе векторов K и K_ν и существует такой вектор смещений $u = (u_x, u_y, u_z)$, что

$$\tau_{xx} = \lambda \text{div } u + 2\mu \epsilon_{xx}, \quad \tau_{xy} = \mu (u_{xy} + u_{yx})$$

и т.п. Здесь λ и μ — постоянные Ляме для рассматриваемой упругой среды.

Пусть \hat{q} — тензор напряжений, действующий для области Ω задачу теории упругости при заданных K и K_ν ; составляющие этого тензора пусть будут $\hat{q}_{xx}, \hat{q}_{xy}, \dots, \hat{q}_{zx}$ и т.д.

что $\hat{q} \in \tilde{H} \ominus H$. Положив $T = V - \hat{q}$, мы сведем вариационную задачу п. 1^о к следующей:

$$\|V - \hat{q}\| = \min, \quad V \in \mathcal{M} = \{V \in H: \|V - \hat{q}\| \leq \mathcal{K}\}; \quad (\text{IX.4.6})$$

$$\|T\| = \tau; .$$

Как обычно, задача (IX.4.6) неразрешима, если $\alpha = \inf_{V \in H} \|V - \hat{q}\| > \mathcal{K}$; она разрешима единственным образом, если $\alpha < \mathcal{K}$; неразрешима или имеет неустойчивое решение, если $\alpha = \mathcal{K}$. Будем считать, что $\alpha < \mathcal{K}$. Тогда существует такой тензор $V_0 \in H$, что $\|q\| < \mathcal{K}$, $\hat{q} = V_0 + \hat{q}$. Обозначим $V = V' - V_0$. Тогда $V - \hat{q} = V' - q$. Далее $\|V - \hat{q}\|^2 = \|V' - q\|^2 = \|V' - V_0 + \hat{q}\|^2$; заменив обозначение V' на V , мы приведем задачу (IX.4.6) к виду

$$\|V - V_0\| = \min, \quad V \in \mathcal{M}; \quad (\text{IX.4.7})$$

$$\mathcal{M} = \{V \in H: \|V - q\| \leq \mathcal{K}\}.$$

Примем допущение, аналогичное тому, которое было принято при анализе той же задачи: тензор упругих напряжений \hat{q} устойчив метриках \mathcal{C} и (IX.4.4) относительно малых возмущений векторов \mathcal{H} и \mathcal{H}_y и области Ω . Тогда (ср. п. 3^о § 2) относительно тех же возмущений устойчив и тензор q . Повторив рассуждения п. 3^о § 2, мы убедимся, что решение трехмерной задачи упруго-пластического состояния устойчиво относительно перечисленных выше малых возмущений, а также относительно малых возмущений постоянной пластичности \mathcal{K} .

3^о. Легко заметить, что общая оценка погрешности аппроксимации (УИ.2.11) верна для задачи настоящего параграфа: если T_* — точное решение этой задачи, а $T_*^{(n)}$ — ее приближенное решение, полученное заменой пространства H его конечномерным подпространством H_n , то

$$\|T_* - T_*^{(n)}\| = O(\nu \epsilon_n(T_*)); \quad \epsilon_n(T_*) = \inf_{T_n \in H_n} \|T_* - T_n\|; \quad (\text{IX.4.3})$$

$$\|T\| = \|T\| + \|T\|.$$

Далее так же верно и общая оценка погрешности выражения (УИ.6.15).

4^о. В заключение отметим некоторые работы, в которых ис-

следуют вариационные неравенства, относящиеся к задачам упруго-пластического состояния. В работе [1] Гаевский рассмотрел две задачи упруго-пластического состояния: о плоском напряженном состоянии и о кручении цилиндрического стержня; обе задачи исследованы в вариационной постановке для случая односвязной области. Каждая из этих задач сведена к уравнению с оператором, удовлетворяющим условию теоремы Браудера о неподвижной точке. Лангенбах [2] получил аналогичный результат для трехмерной задачи. Сценки устойчивости для решений некоторых вариационных неравенств дал Гаевский [3]; в качестве примера им рассмотрено, в частности, прямолинейное течение вязко-пластической жидкости в трубе. В работе [4] того же автора принцип Хаара-Кармана используется в рамках теории пластичности Прандтля-Рейсса.

ЛИТЕРАТУРА

Агмон С., Дуглио А., Шренберг Л. И. Оценки решений эллиптических уравнений вблизи границы. Перев. с англ. 1962, М., с.205.

Атакишиева Р.Х. И. Устойчивость метода Бунцова-Галеркина для уравнения с вполне непрерывным оператором в банаховом пространстве. Докл.АН Азерб.ССР, 1966, 2, № 1, 7-11.

Бабушка И., Витесек П., Прагер М. И. Численные процессы решения дифференциальных уравнений. Перев. с англ. 1969, М., с. 368.

Балакришнан А.В. I Прикладной функциональный анализ. Перев. с англ. 1980, М., с.383.

Барн Н.К. I. Обобщение неравенств С. Липштейна и А.А.Маркова. Изв. АН СССР, серия матем., 1961, 18, № 2, 159-176.

Березин Л.С., Кидков Н.И. I. Метод вычислений, т.1, 1966, М., с.631

Бирман М.Л. I. О методе Фридрикса расширения положительно определенного оператора до сопряженного. Зап.Ленингр.Горн. ин-та, 1956, 33, № 3, 132-136.

Вазов В., Форсайт Дж.И. Разностные методы решения дифференциальных уравнений в частных производных. Перев. с англ., 1963, М., с.487.

Вайникко Г.М. I. Оценки погрешности метода Галеркина для линейного дифференциального уравнения. Уч.зап.Тартуск. ун-та, 1962, вып. 129, Труды по матем. и мех., III, 394-416.

--- 2. Оценка погрешности метода Рунге для линейного однородного уравнения. Уч. зап. Тартуск. ун-та, 1962, вып.129. Труды по матем. и мех., III, 417-427.

--- 3. Некоторые оценки погрешности метода Бунцова-Галеркина. I. Асимптотические оценки. Уч. зап. Тартуск. ун-та, 1962, вып. 150. Труды по матем. и мех., IV, 188-201.

--- 4. Некоторые оценки метода Бушова-Галеркина. II. Оценка n -го порядка сходимости. Уч. зап. Тартуск. ун-та, 1964, вып.150. Труды по матем. и мех., I., с. 2-5.

--- 5. Асимптотические оценки погрешности проекционных методов в проблеме сглаженных значений. ВМММ, 1964, 4, № 405-5.

- - - 6. Оценки погрешности метода Бубнова-Галеркина в проблеме собственных значений. ЖРММЖ, 1965, 5, № 4, 587-607.

- - - 7. Необходимые и достаточные условия устойчивости метода Галеркина - Петрова. Уч. зап. Тартуск. ун-та, 1965, вып. 177, 141-147.

- - - 8. О сходимости и устойчивости метода коллокации. Диф. ур-ия, 1965, 1, № 2, 244-254.

- - - 9. О быстроте сходимости некоторых приближенных методов типа Бубнова - Галеркина в проблеме собственных значений. Изв. ВУЗ. Математика, 1966, № 2, 35-45.

- - - 10. О сходных операторах. ДАН СССР, 1968, 179, № 5, 1029-1031.

- - - 11. Анализ дискретизационных методов. Тарту, 1976, с. 181.

Вайникко Г.М., Михлен С.Г. I. Резольвента Фредгольма и обращение матрицы, линейно зависящей от параметра. Уч. зап. Тартуск. ун-та, 1978, вып. 448. Труды по матем. и мех., XXI, 94-98.

Варга Р.С. I. Функциональный анализ и теория аппроксимации в численном анализе. Перев. с англ. 1974, М., с. 126.

Велиев М.А. I. Исследования устойчивости метода Бубнова - Галеркина для нестационарных задач. ДАН СССР 1964, 157, № 1, 16-19.

- - - 2. Об устойчивости метода Бубнова - Галеркина для нестационарных задач. В сборн. "Вопросы вычисл. матем. и вычисл. техн.", тзд. АН Аз. рб. ССР, 1964, вып. 3.

- - - 3. Об устойчивости метода Бубнова - Галеркина для уравнения n -го порядка в гильбертовом пространстве. Сб. матем. ж., 1968, IX, № 3, 783-789.

- - - 4. Об устойчивости метода Бубнова - Галеркина для уравнений первого порядка с переменными коэффициентами в гильбертовом пространстве. Диф. ур-ия, 1969, 5, № 3, 473-487.

- - - 5. Некоторые достаточные условия устойчивости метода Бубнова - Галеркина для уравнений второго порядка в гильбертовом пространстве. Уч. зап. Азерб. ун-та, серия физ.-мат. наук, 1972, № 2, 17-24.

- - - 6. К устойчивости метода Бубнова - Галеркина для одного класса аппроксимируемых параболических уравнений. Научн. труды

Азерб. ун-та; Вопросы прикл. матем. и киберн., № 1979, № 1, 30-34.

- - - 7. Об устойчивости метода конечных элементов для параболических уравнений. Докл. АН Азерб.ССР, 1981, 37, № 4, 3-6.

- - - 8. К устойчивости метода Бубнова - Галеркина для линейных уравнений в гильбертовом пространстве. Докл. АН Азерб. ССР, 1981, 37, № 5, 3-6.

- - - 9. Устойчивость метода Бубнова - Галеркина для линейных уравнений второго порядка с переменными коэффициентами. Докл. АН Азерб.ССР, 1981, 37, № 7, 8-11.

Воеводин В.В. I. Ошибки округления и устойчивость в прямых методах линейной алгебры. 1969, М., с.153.

Вулик Б.З. I. Краткий курс теории функций вещественной переменной. 1973, М., с.350.

Габдулхаев Б.Г. I. Об одном прямом методе решения интегральных уравнений. Изв. ВУЗ, Математика, 1965, № 3, 51-60.

- - - 2. Приближенное решение сингулярных интегральных уравнений методом механических квадратур. ДАН СССР, 1968, 179, № 2, 260-263.

- - - 3. Оптимальные аппроксимации решений линейных задач. 1980, Казань, с.231.

Габдулхаев Б.Е., Душков П.Н. I. Метод механических квадратур для сингулярных интегральных уравнений. Изв. ВУЗ, Математика, 1974, № 12 (151), 3-14.

Гавурий М.К. I. Лекции по методам вычислений. 1971, М., с.248.

Галин Л.А. Упруго-пластическое кручение стержней полигонального сечения. Прикл. матем. и мех., 1944, 8, № 4, 307-322.

- - - 2. Упруго-пластическое кручение призматических стержней. Прикл. матем. и мех., 1949, 13, № 3, 287-296.

Гельман И.В.К задаче о минимуме нелинейного функционала. Уч. зап. Лен. Пед. ин-та им. Герцена, 1958, 166, 255-263.

Гловински Р., Лиоис Ж.-Л., Тремольтер Р. I. Численное исследование вариационных неравенств. Перев. с франц., 1979, М., с.574.

Голунгов С.К., Рябенский В.С. I. Введение в теорию разностных схем. 1962, М., с.340.

Головкин К.К. I. О приближении функций в произвольных нормах. Труды матем. ин-та Л.М.Стеклова АН СССР, 1964, 70, 26-37

Гохберг И.Ц., Фельдман И.А. I. Уравнения в овертках и проекционные методы их решения. 1971, М., с.352.

Гурса Э. I. Курс математического анализа, т.Ш, ч.П. Перев. с франц., 1934, М., с.318.

Гус: ан Ю.А., Оганезян Л. A. I. Оценка сходимости конечноразностных схем для вырожденных эллиптических уравнений. ЖВММФ, 1966, 5, № 2, 351-357.

Демьянович Ю.Г. I. Устойчивость метода сезок для эллиптических задач. ДАН СССР, 1965, 134, № I, 20-23.

- - - 2. Об устойчивости проекционных методов. Зап.научн.семина.ЛОМИ АН СССР, 1971, 23, 5-15.

- - - 3. Об устойчивости и длительности вычислений в вариационно-разностном методе. Зап.научн.семина.ЛОМИ АН СССР, 1978, 80, 5-29.

- - - 4. Об оценках окростности сходимости проекционных методов решения некоторых вариационных неравенств. Деп. в ВИНИТИ, 1980, № 1247-80 ДЕП, с.47

Демьянович Ю.К., Михлин С.Г. I. О осточной аппроксимации функций соболевских пространств. Зап.научн.семина.ЛОМИ АН СССР, 1973, 35, 6-II.

Джилшариани А.В. I. Некоторые оценки погрешности метода Рунца - Галеркина для обыкновенных дифференциальных уравнений. Труды Тбил. матем. ин-та, 1959, 25, 285-306.

- - - 2. Оценка погрешности метода Рунца для неоднородного дифференциального уравнения. Сообщ. АН Груз.ССР, 1960, 25, № 3, 257-272.

- - - 3. О быстроте сходимости приближенного метода Рунца. ЖВММФ, 1963, т.3, № 4, 354-360.

- - - 4. О быстроте сходимости метода Бубнова - Галеркина. ЖВММФ, 1964, 4, № 2, 343-348.

- - - 5. О быстроте сходимости обобщенных методов Рунца - Бубнова - Галеркина. Сообщ. АН Груз. ССР, 1967, 40, № 3, 553-560.

- - - 6. О методах Рунца и Бубнова - Галеркина. Сообщ. АН Груз.ССР, 1967, 40, № I, 11-16.

- - - 7. О методе Бубнова - Галеркина. ЖВМиМФ, 1967, 7, № 6, 1398-1402.

- - - 8. О методе Рунца для одного нелинейного уравнения. Сосощ. АН Груз. ССР, 1968, 51, № 1, 19-24.

- - - 9. О методах наименьших квадратов и Бубнова - Галеркина. ЖВМиМФ, 1968, 8, № 5, IIII-III6.

- - - 10. О приближенном методе Рунца для одного нелинейного уравнения. Труды Тбил. матем. ин-та, 1969, 36, 29-46.

- - - 11. Устойчивость метода Рунца для одного нелинейного уравнения. ЖВМиМФ, 1970, 10, № 4, 841-847.

- - - 12. К решению сингулярных интегральных уравнений приближенными проекционными методами. ЖВМиМФ, 1979, 19, № 5, II49-II61.

- - - 13. О сходимости невязки в методе конечных элементов. ДАН СССР, 1980, 254, № 5, 1052-1055.

- - - 14. К решению сингулярных интегральных уравнений коллокационными методами. ЖВМиМФ, 1981, 21, № 2, 355-362.

Дочбѣт Д.Н. I. Устойчивость метода Рунца для задач спектральной теории операторов. Труды матем. ин-та им.Стеклова АН СССР, 1965, 84, 78-92.

Дюво Г., Дюво Ж.-Л. I. Неравенства в механике и физике. Перев. с франц. 1980. М., с.383.

Запрудский Я.М. I. Свойства метода Гальборкина - Петрова у банаховому простору. Докл.АН УССР, серия А, 1972, № 4, 309-312.

Ильин В.П. I. О сходимости вариационных процессов. ДАН СССР, 1951, 81, № 2, 137-140.

- - - 2. Оценка погрешности в методе Рунца для обыкновенных дифференциальных уравнений. Труды матем. ин-та им.Стеклова АН СССР, 1959, 53, 43-63.

- - - 3. Некоторые неравенства в функциональных пространствах и их применение к исследованию сходимости вариационных процессов. Труды матем. ин-та им.Стеклова АН СССР, 1959, 53, 64-127.

- - - 4. Свойства некоторых классов дифференцируемых функций многих переменных, заданных в n -мерной области. Труды матем. ин-та им.Стеклова АН СССР, 1962, 66, 227-363.

- - - 5. Некоторые замечания о сходимости последователь-

ности функций п-линоммального типа в пространствах $W_p^{(k)}(G)$. Зап. научн.семина.ЛОМИ АН СССР, 1968, II, 73-96.

Канторович Л.В. I. Об одном методе приближенного решения дифференциальных уравнений в частных производных. ДАН СССР, 1934, 2, 532-536.

- - - 2. Об одном эффективном методе решения экстремальных задач для квадратичных функционалов. ДАН СССР, 1945, 48, № 7, 455-460.

- - - 3. Функциональный анализ и прикладная математика. УМН, 1948, 3, № 6(28), 89-185.

Канторович Л.В., Акилов Г.П. I. Функциональный анализ, изд. 2-е, М., 1977. с.741.

Канторович Л.В., Крылов В.И. I. Приближенные методы высшего анализа, изд. 3-е, Л.-М., 1949 с. 695.

Каршлюковская Э.Б. I. О сходимости интерполяционного метода для обыкновенных дифференциальных уравнений. УМН, 1953, 8, № 3 (55), III-II8.

- - - 2. О сходимости метода коллокации. ДАН СССР, 1963, 151, № 4, 736-769.

- - - 3. О сходимости метода коллокации для некоторых граничных задач математической физики. Сиб.матем.ж., 1963, 4, № 3, 632-640.

Колмогоров А.Н., Фомин С.В. I. Элементы теории функций и функционального анализа. Изд. 5-е. ГИИ, М., с. 542.

Коряев В.Г. I. Некоторые вопросы построения и исследования схем метода конечных элементов. Числ.методы мех. сплошн. среды, 1974, № 5, 59-87.

Корнеев В.Г., Пономарев С.Э. I. Применение криволинейных конечных элементов в схеме решения линейных алгебраических уравнений порядка $2n$. Числ.методы мех. сплошн. среды, 1974, 5, № 5

Красносельский М.А. I. Сходимость метода Галеркина для нелинейных уравнений. ДАН СССР, 1950, 73, № 6, II21-II24.

Красносельский М.А. и др. I. Приближенное решение операторных уравнений. ГИИ, М., с. 455.

Крылов В.И. I. Избранные труды. т.3, 1961, Киев, с.398.

Курцадзе И.И. I. Трехмерные задачи математической теории упругости и термоупругости. 1971, М., с.333.

Ладженская О.А., Уральцева Н.Н. I. Линейные и квазилинейные уравнения эллиптического типа, изд. 2-е, 1975, М., с.573.

Лионо Ж.-Л. I. Оптимальное управление системами, описываемыми уравнениями в частных производных. Перев. с франц., 1972, М., с.474.

- - - 2. Некоторые методы решения нелинейных краевых задач. Перев. с франц., 1972, М., с. 587.

Мамагов А.З. I. Применение метода Галеркина к некоторому квазилинейному уравнению параболического типа. Вестн. ЛГУ, 1981, № 13, 37-45.

- - - 2. О погрешности и устойчивости метода Бубнова - Галеркина для квазилинейных параболических задач с производной по времени в граничном условии. Док. в ВИН.ПИ, 1981, № 146-81 ДСП, с.36.

Марчук Г.И. I. Методы вычислительной математики, 1973, Новосибирск, с. 352.

Митлин С.Г. I. Интегральные уравнения и их приложения, М.-Л., 1949, с. 80.

- - - 2. Проблема минимума квадратичного функционала. 1952, М., с.213.

- - - 3. Вариационные методы в математической физике, изд. 2-е, 1970, М., с.512.

- - - 4. Об устойчивости метода Рунге. ДАН СССР, 1960, 135, № 1, 16-19.

- - - 5. Некоторые условия устойчивости метода Рунге. Вестн. ЛГУ, 1961, № 14, 40-51.

- - - 6. Многомерные сингулярные интегралы и интегральные уравнения. 1962, М., с. 254.

- - - 7. Об устойчивости некоторых вычислительных процессов. ДАН СССР, 1964, 157, № 2, 271-273.

- - - 8. Численная реализация вариационных методов. 1966, М., с. 432.

- - - 9. О функции Кессера. В сборн.: "Проблемы матем. анализа", изд. ЛГУ, 1966, 39-49.

- - - 10. Курс математической физики. 1968, М., с. 575.

- - - 11. Неравенства типа А.А.Маркова и полиномиальные приближения к Рунге. В книге: "Дифференциальные уравнения с частными производными", 1970, М., 153-169.

- - - 12. О вариационно-разностном методе для многосмерных краевых задач. Зап. научн. семинар. ЛОМИ АН СССР, 1971, 23, 99-114.
- - - 13. О вариационно-разностном методе для одномерных краевых задач. ДАН СССР, 1971, 198, в 1, 39-41.
- - - 14. О сеточной аппроксимации вырождающихся одномерных дифференциальных уравнений второго порядка. Вестн. ЛГУ, 1973, в 1, 52-67.
- - - 15. Спектр пучка операторов теории упругости. УМН, 1973, 28, в 5(171), 43-82.
- - - 16. О числе обусловленности вариационно-сеточной матрицы. Вестн. ЛГУ, 1973, в 13, 167-164.
- - - 17. Вариационно-сеточная аппроксимация. Зап. научн. семинар. ЛОМИ АН СССР, 1974, 48, 32-188.
- - - 18. Об одном методе приближенного решения интегральных уравнений. Вестн. ЛГУ, 1974, в 13, 26-33.
- - - 19. О новом методе приближенного решения интегральных уравнений. Вестн. ЛГУ, 1976, в 1, 35-44.
- - - 20. Замечания о резольвентном методе. Зап. научн. семинар. ЛОМИ АН СССР, 1976, в 56, 5-13.
- - - 21. О приближенном решении интегральных уравнений со слабой особенностью. Вестн. ЛГУ, 1976, в 13, 31-36.
- - - 23. Об одном варианте резольвентного метода. Вестн. ЛГУ, 1978, в 1, 54-59. — ? 22
- - - 24. О постоянных множителях в оценках погрешности вариационно-сеточной аппроксимации. Зап. научн. семинар. ЛОМИ АН СССР, 1978, 80, 125-166.
- - - 25. О методе наименьших квадратов для многосмерных сингулярных интегральных уравнений. В сборн. "Комплексный анализ и его прилож.". К 70-летию акад. И. Н. Векуа, 1978, м., 401-408.
- - - 26. О приближенном решении односторонних вариационных задач. УМН, 1979, 34, в 4(206), с. 166.
- - - 27. О приближенном решении односторонних вариационных задач. Изв. ЛГУ, Математика, 1980, в 2(213), 45-48.
- - - 28. О погрешности приближенного решения односторонних вариационных задач. Вестн. ЛГУ, 1980, в 15, 46-50.
- - - 29. Об устойчивости решений односторонних вариационных задач; приложения к теории пластичности. Зап. научн. семинар.

Мин. ДОМЫ АН СССР, 1980, 102, 68-101.

- - - 30. погрешность искажения в простейшей* односторонней вариационной задаче. Вестн. ЛГУ, 1981, № 13, 45-50.

- - - 32. О погрешности метода Бубнова - Галеркина для аналитических краевых задач. Зап. науч. семина. ДОМЫ АН СССР, 1981, III, 137-144.

- - - 33. О погрешностях вычислительных процессов I. Изв. ВУЗ. Математика, 1981, № 7(230), 62-71.

- - - 34. О погрешностях вычислительных процессов. II. Изв. ВУЗ. Математика, 1981, № 8(231), 32-38.

Михлин С.Г., Радова Р.К., I. О приближенном решении сингулярных интегральных уравнений. Изв. ВУЗ, Математика, 1974, № 5 (144). 158-162.

Михлин С.Г., Смолицкий Х.Г. I. Приближенные методы решения дифференциальных и интегральных уравнений. 1965, М., с.388.

Мухомелишвили Н.И. I. Некоторые задачи теории упругости, тед. 3-е. 1940, М., с.335.

Надаи А. I. Пластичность. Перев. с англ., 1936, М.-Л., с.280.

Натансон И.П. I. Конструктивная теория функций. 1949, М.-Л., с.688.

Осан Я.-П. I. Приближенное решение эллиптических краевых задач. 1977, М., с.383.

Отгнесян Л.А., Руховец Л.А. I. Вариационно-разностные методы решения эллиптических уравнений, 1979, Ереван, с.335.

Оден Дж. I. Конечные элементы в нелинейной механике сплошных сред. Перев. с англ., 1976, М., с.464.

Перлин П.И. I. Упруго-пластическое кручение стержней овального поперечного сечения. Инж. сборник, 1961, XI, 202-205.

Руховец Л.А. I. К вопросу о построении вариационно-разностных схем для эллиптических уравнений. ВММ, 1972, 13, № 3, 795-798.

Смирнов В.И. I. Курс высшей математики, т. III, ч. I. 1949, М.-Л., с.335.

Соболев С.Г. 2. Введение в теорию кубатурных формул. 1974, М., с.308.

Странг Г. Фитс Дж. 2. Теория методов конечных элементов. Перев. с англ., 1977, М., с.349.

Суворов С.Г. I. К вопросу об устойчивости вычислительного процесса. Вестн. ЛГУ, 1972, № 7, 115-116.

Сьарле Ф. I. Метод конечных элементов для эллиптических задач. Перев. с англ., 1980, М., с. 512.

Талл кин А.Т. I. Системы элементов гальбертова пространства и ряды по ним. Матем.оборн., 1951, 29(71), № 1, 79-120.

Тесия пластичнос я. Сборник статей. I. Перев. с англ., фр. и нем., 1948, М., с. 452.

Тополянский Д.Б., Запрудский Я.М. I. Исследование устойчивости метода Бубнова-Галеркина для некоторых операторных уравнений с переменными коэффициентами. Укр.матем.ж., 1974, 26, №5, 621-634.

Тыттин В.Б. I. Об оценке погрешности приближенных решений односторонних вариационных задач. Вестн. ЛГУ, 1981, № 13, 65-69.

- - - 2. О скорости сходимости приближенных методов решений односторонних вариационных задач. I. Вестн. ЛГУ, 1982, № 13, III-III3.

Уилкинсон Дж.Х. 2. Алгебраическая проблема собственных значений. Перев. с англ., 1970, М., с. 564.

Фаддеев Д.К., Фаддеева В.Н. I. Вычислительные методы линейной алгебры, изд. 2-е, 1967, М.-Л., с. 734.

Харрик И.Ю. I. О приближении функций, обращающихся на границе в нуль, функциями особого вида. Матем. сборн., 1955, 37, 353-384.

- - - 2. О приближении функций, обращающихся в нуль вместе с градиентом на границе области, функциями особого вида. Матем. Сборн., 1959, 47(89), № 2, 177-208.

- - - 3. О приближении функций, обращающихся в нуль на границе области вместе с частными производными, функциями особого вида. Сиб.матем.ж., 1968, 4, № 2, 408-425.

Черноусько Ф.Д., Баничук Н.Л. I. Вариационные задачи механики и упругости. 1973, М., с. 238.

Чистот К.Е. I. Вариационно-сеточная аппроксимация на регулярной треугольной сетке. Изв. ВУЗ, Математика, I 81, № II, 58-63.

- - - 2. Устойчивость вариационно-сеточного процесса на симплектических сетках. Изв. ВУЗ, Математика, 1982 № 4, 51-65.

Шапошникова Т.О. I. Некоторые оценки сходимости вариационных методов. Вестн. ЛГУ, 1971, № 19, 64-69.

- - - 2. Априорные оценки погрешности некоторых вариационных методов. Диссертация, ЛГУ, 1973, с.93.

- - - 3. Априорные оценки погрешности некоторых вариационных методов в соболевских пространствах. В сборн.: Вариационно-разностные методы в матем. физ. СО АН СССР, 1976, Новосибирск, 182-194.

- - - 4. Априорные оценки погрешности вариационных методов в банаховых пространствах. ЖВМиМФ, 1977, 17, № 5, 1144-1152.

Искова Е.Н., Яковлева М.Н. I. Некоторые условия устойчивости метода Петрова - Галеркина. Труды Матем. ин-та им.Стеклова АН СССР, 1962, 66, 182-189.

Agmon S., Douglis A., Nirenberg L. 2. Estimates near the boundary for solutions of elliptic partial differential equations satisfying general boundary conditions. Comm. on pure and appl. Math., 1964, XVII, 35-92.

Andersen R.S., de Hoog G.R., Lucas M.A. (ed.). 1. The application and numerical solution of integral equations. Alphen aan den Rijn, 1980, pp.259.

Anselon 1. Collectively compact operator approximation theory and applications to integral equations. Englewood Cliffs, 1971, pp.138.

Berger G., Mox M., 1. Interpolation über ein Dreiecksnetz. Wiss. Zeitschr. d.Univ. Halle, 1976, 25, N 1, 43-48.

Bertram G. 1. Verschärfung einer Fehlerabschätzung zum Ritz-Galerkinschen Verfahren von Kryloff für Randwertaufgaben. Numerische Math., 1959, 1, N 3, 135-141.

Börsohn-Supan W. 1. Bemerkungen zur Fehlerabschätzung beim Ritz - Galerkin Verfahren nach Krylow Numerische Math., 1960, 2, N 2, 79-83.

Bramble J.H., Zlamal M. 1. Triangle elements in the finite element method. Math. of Comput., 1970, 21, 1, 112, 809-820.

Brezis H. 1. Problèmes unilatéraux. J. de Math. pure et appliquée. 1972, 51, 1-168.

Brezis H., Stampacchia G. 1. Sur la régularité de la solution d'inéquation elliptiques. Bull.Soc.Math. de France, 1968, 96, 153-180.

Brezzi F., Nager W.W., Raviart P.A. 1. Error estimates for the finite element solution of variational inequalities. Part I, Primal theory, *Numerische Math.*, 1977, 28, 431-443.

Calderon A.P. 1. Lebesgue spaces of differentiable functions and distributions. *Partial differential equations Amer. Math. Soc. Proc. of symposia in pure mathematics*, 1961, IV, 33-49

Courant R. 1. Variational methods for the solution of equilibrium and vibrations. *Bull. Amer. Math. Soc.*, 1943, 49, N 1, 1-23.

Delves L.M., Walsh J. (ed.). 1. Numerical solution of integral equations. Oxford, 1974, pp. 333.

El Kholli A. 1. Régularité de la solution et majoration de l'erreur d'approximation pour un problème de Dirichlet non linéaire. *Compt. Rend Acad. Sci. Franç. Ser A*, 1973, 277, 735-737.

Elliott C.M. 1. On the finite element approximation of an elliptic variational inequality arising from an implicit time discretisation of the Stefan problem. *IMA Journ. of numer. analysis*, 1981, 1, N 1, 115-126.

Falk R.S. 1. Error estimates for the approximation of a class of variational inequalities. *Math. Comput.*, 1974, 28, 963-971.

- - - 2. Approximation of an elliptic boundary value problem with unilateral constraints. *RAIRO, Anal. Numer.*, 1975, 9, 5-12.

Gajewski H. 1. Zur Lösung einer Klasse konvexer Minimalprobleme der Plastizitätstheorie. *Zeitschr. f. angew. Math. und Mech.*, 1969, 49, N 1/2, 83-89.

- - - 2. Zur numerischen Stabilität des Ritzschen Verfahrens bei nichtlinearen Gleichungen. *Math. Nachr.*, 1970, 43, N 1-6, 261-311

- - - 3. Stabilitätsaussagen für einige ungelöste Randwertprobleme. *Zeitschr. f. angew. Math. und Mech.*, 1977, 57, 439-447.

- - - 4. Ein Verfahren zur Ermittlung elastisch-plastischer Spannungsfelder. *Beit. z. Spannungs- und Dehnungsanalyse*, herausg. von K. Schröder, Berlin, 1968, V 15-32.

Gołwinski P., Larchon J. 1. Torsion élasto-plastique d'une barre cylindrique de section multi-connexe. *Journ. de mécanique*, 1973, 12, 151-171.

Glowinski R., Marocco A. 1. Sur l'approximation par éléments finis d'ordre un, et la relation par pénalisation-dualité d'une classe de problèmes de Dirichlet non linéaires. RAIRO, Anal.numér., R-2, 1975, 41-76.

Golub G.M. 1. Bounds for round-off errors in the Richardson second order method. Nordisk Tidskrift for Informationsbehandling (BIT), 1962, 2, H. 4, 212-223.

Haar A., v.-Karmen T. 1, Zur Theorie der Spannungszustände in plastischen und sandartigen Zuständen. Nachr. d. Königlich. Gesellsch. d. Wissensch. zu Göttingen, 1909, H. 2, 204-218.

Hencky H. 1. Über einige statisch bestimmte Fälle des Gleichgewichts in plastischen Körper. Zeitschr. f. angew. Math. und Mech., 1921, 3, H. 4, 241-251.

Hille E., "zeg" G., Tamarkin J.D. 1. On some generalizations of a theorem of A.Markoff. Duke Math.J., 3, 1937, 729-739.

Lanchon H. 1. Solution du problème de torsion d'une barre cylindrique de section quelconque. Comptes Rend. de l'Acad. d. Sci.Franc, ser.A, 1969, 269, 791-794.

Langenbach A. 1. Die Linearisierung nichtlinearen Gleichungen. Math.Nachr., 1962, 24, 33-51.

- - - 2. Monotone Potentialoperatoren. 1976, Berlin, SS. 358.

Lehman N.J. 1. Eine Fehlerabschätzung zum Ritzschen Verfahren für inhomogene Randwertaufgaben. Numerische Math., 1960, 2, N 2, 60-66.

Lévy M. 1. Mémoire sur les équations générales des mouvements intérieurs des corps solides ductiles au delà des limites ou l'élasticité pourrait les ramener à leur premier état. Comptes Rend. de l'Acad. d.Sci.Franc., 1870, 70, 13.3-1325.

Michlin S.G. 22 Approximation auf dem kubischen Gitter. Berlin, 1976, SS.197

- - - 31. Konstanten in einige Ungleichungen der Analysis. Leipzig, 1981, SS.119.

Michlin S.G., Prädorf S. 1. Singuläre Integraloperatoren. Berlin 1980. SS.514.

v.-Mises R. 1. Mechanik der festen Körper in plastisch-deformablen Zustand. Nachr. v. d. Königlich. Gesellsch. d.

Wissensch. zu Göttingen, 1913, N. 4, S.582.

- - - 2. Bemerkungen zur Formulierung des mathematischen Problems der Plastizitätstheorie. *Zeitschr. f. angew.Math. und Mech.*, 1925, IV, H. 2, 147-149.

Mittelman H.D. 1. On the approximate solution of nonlinear variational inequalities. *Numerische. Math.*, 1978, 29, 451-462.

Morrey B.Jr. 1. Multiple integrals in the calculus of variation N.Y., 1966, pp.506.

Mosco U. 1. Implicit variational problems and quasi variational inequalities. *Lecture Notes in Math.*, 1976, 543, 83-136.

Nitsche J. 1. L-convergence of finite element approximation. *Lecture Notes in Math.*, 1977, 606, 261-274.

Omodei B.J. 1. Stability of the Rayleigh - Ritz - Galerkin procedure for elliptic boundary value problem. A thesis submitted to the Australian National Univ. 1975, pp.138.

Proßdorf S., Schmidt G. 1. A finite element collocation method for singular integral equations. *J. Math.Machr.*, 1981, 100, 35-60.

Proßdorf S., Silberman B. 1. Projektionsverfahren und die näherungsweise Lösung singulärer Gleichungen. Leipzig, 1977, SS.226.

Ritz W. 1. Über eine neue Methode zur Lösung gewisser Variationsprobleme der Mathematischen Physik. *J. f. d. reiner und angew. Math.*, 1908, 135, H.1.

de-Saint-Venant B. 1. Sur l'établissement des équations des mouvements intérieurs opérés dans les corps ductiles au delà des limites où l'élasticité pourrait les ramener à leur premier état. *Comptes Rend. de l'Acad. d. Sci. Franq.*, 1870, 70, 473-480.

Scarpini F. 1. Some nonlinear complementary systems algorithms and applications to unilateral boundary value problems. *Ann.Math.pura ad appl.* 1980, 123, 185-202.

Scarpini F., Vivaldi M.A. Error estimates for the approximation of some unilateral problem. *RAIRO, Anal.numer*, 1977, 11, N 2, 197-208.

Sobolev S. 1. The problem of propagation of plastic state. *Труды Сейсмологического института АН СССР*, № 49, 1935, pp. 15

Strang G. 1. Approximation in the finite element method

Numerisch Math., 1972, 19, 91-98.

Strang G., *ix* 3. 1. A Fourier analysis of the finite element method. (Препринт; год издания не указан).

Tucker T.S. Stability of nonlinear computing schemes. *SIAM J. Numer. anal.*, 1969, 6, 72-81.

Temple G. The general theory of relaxation methods applied to linear systems. *Proc. of Royal Soc., ser. A*, 1939, 16^o, 476-500.

Wilkinson J.H. 1. Rounding errors in algebraic processes. *Nat. Phys. Laboratory. Notes on appl. science. Englewood Cliffe*, 1963, pp. 161.

Woźniakowski H. 1. Round off error analysis of iterations for large linear systems. *Numerische Math.*, 1978, 30, fasc. 3, 301-314.

Zenisek A. 1. Interpolation polynomials on the triangle. *Numerische Math.*, 1970, 15, 283

- - - 2. Hermite interpolation on simplexes in the finite element method. *Proc. EQUADIFF-III, Rome 1972*, p. 271.

- - - 3. Polynomial approximations on tetrahedrons in the finite element method. *J. Approx. Theory*, 1973, 7, N 4, 334-351.

Zienkiewicz M. 1. On the finite element method. *Numerische Math.* 1968, 1^o, 394-408.

S.G. Mikhlín

Errors of numerical processes

Summary

This book deals with a sufficiently wide class of problems of the numerical mathematics and the numerical processes leading to the construction of the solution of a given problem. The sources of the errors related with these processes are investigated under the assumption that the given problem is posed exactly. The author thinks that in many cases there are exactly four sources of errors: approximation error, resulting from the change of the given problem by a less simplified problem; distortion error which is connected with the fact that the data of the simplified problem are calculated with some errors; algorithm error, which arise in the cases when the algorithm, used for solve the simplified problem, gives only an approximate solution after a finite number of steps; round-off error, which arise because of the errors of the fulfilment of the instructions of the mentioned algorithm. Estimates of all these errors are given for a wide class of numerical processes.

The book contains 9 chapters. The Chapter I contains the statement of the problem and the investigation of two simpler examples. Errors of linear numerical processes are investigated in the Chapt. II-VI and the errors of non-linear ones are considered in the Chapt. VII-IX.

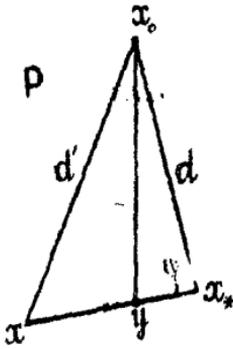


Рис. 1

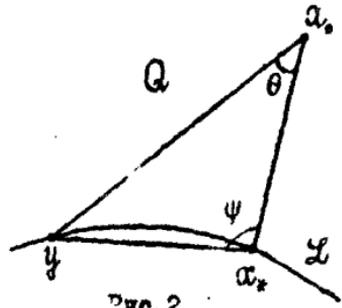


Рис. 2

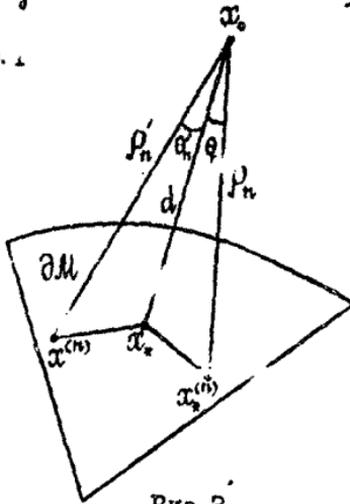


Рис. 3

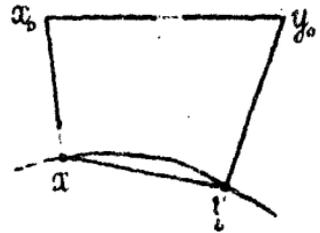


Рис. 4

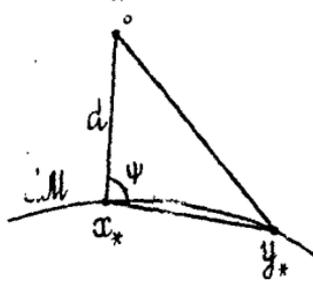


Рис. 5

ОГЛАВЛЕНИЕ

Предисловие	3
Глава I. Постановка проблемы	7
§ 1. Задачи вычислительной математики. Вычислительные процессы	7
§ 2. Источники погрешностей	11
§ 3. Погрешности квадратурных формул	14
§ 4. Погрешности решения систем линейных алгебраических уравнений	17
Глава II. Погрешность аппроксимации.	23
§ 1. Метод Рунге. Оценка в энергетической метрике	23
§ 2. Некоторые обобщения теоремы Джексона на функции многих переменных	26
§ 3. Обыкновенные дифференциальные уравнения	28
§ 4. Уравнения эллиптического типа. Оценки в энергетической норме	30
§ 5. Некоторые оценки погрешности производных. Литературные ссылки	31
§ 6. Некоторые марковские неравенства	32
§ 7. Оценки аппроксимации старших производных	36
§ 8. Оценки аппроксимации младших производных	42
§ 9. Метод Бунцова - Галеркина	47
§ I0. Разностные методы	49
§ II. Метод коллокации	51
Глава III. Погрешность искажения	56
§ 1. Погрешность искажения свободного вычислительного процесса	56
§ 2. Устойчивость свободного процесса относительно искажения	57
§ 3. Устойчивость процесса Рунге	62
Примеры	68
§ 4. Об устойчивости процесса Бунцова - Галеркина	76
Погрешность искажений рекуррентного численного процесса	81
Устойчивость рекуррентного вычислительного процесса	82
Устойчивость искажения и устойчивость метода Гаусса	89
Устойчивость метода Рунге	90

Глава IV. Погрешности алгоритма и округления	95
§ 1. Число обусловленности	95
§ 2. Погрешность алгоритма в итерационном процессе	96
§ 3. Погрешность округления в рекуррентном процессе	98
§ 4. Погрешность округления метода наискорейшего спуска	101
§ 5. Погрешность округления метода Рундсона для линейных или обратных систем	104
Глава V. Погрешности метода конечных элементов	107
§ 1. Обзор некоторых результатов: погрешность аппроксимации	107
§ 2. Обзор некоторых результатов: погрешность искажения и число обусловленности и тогда конечных элементов	112
§ 3. Метод конечных элементов с одной кусочно полиномиальной исходной функцией. Постоянный множитель в оценке погрешности аппроксимации	115
§ 4. Погрешность искажения	119
§ 5. Симплициальные сетки. Оценка аппроксимации	122
§ 6. Теорема об устойчивости	125
Глава VI. Погрешности приближенного решения интеграль- ных уравнений	128
§ 1. Погрешности метода механических квадратур для уравне- ний Фредгольма	129
§ 2. Уравнения Фредгольма, решаемые итерациями	133
§ 3. Погрешности определения резольвенты Фредгольма	137
§ 4. Сведения к алгебраической системе	142
§ 5. Сингулярные интегральные уравнения. Метод наименьших квадратов	145
§ 6. Методы механических квадратур и комбинации для одно- мерных сингулярных уравнений	148
Глава VII. Нелинейные вычислительные процессы	152
§ 1. О погрешности аппроксимации метода Рундса	152
§ 2. Погрешность искажения и устойчивость свободного нели- нейного процесса	156
§ 3. Об устойчивости процессов Рундса и конечных элементов	159
§ 4. Погрешности дискретизации и округления рекуррентного вы- числительного процесса	166
§ 5. Метод Ньютона - Канторовича	17

Глава VIII. Некоторые односторонние вариационные задачи	174
§ 1. Постановка задачи и ее приближенное решение	174
§ 2. Погрешность аппроксимации	179
§ 3. Случай, когда оценка погрешности аппроксимации улучшается	182
§ 4. Оценка погрешности аппроксимации для более общей задачи	190
§ 5. Об устойчивости точного решения односторонней вариационной задачи	192
§ 6. Оценка погрешности искажения	201
Добавление к главе VIII (Б. В. Тихтин)	205
Глава IX. Упруго-пластическое состояние по Сен-Венану-Мизесу и Хаару-Карману	214
§ 1. Уравнения Сен-Венана и Мизеса. Вариационный принцип Хаара - Кармана	214
§ 2. Кручение стержня	227
§ 3. Плоская задача	230
§ 4. Трехмерная задача	236
Цитированная литература	240

УДК 519.6

Погрешности вычислительных процессов. Михлин С.Г.
Институт прикладной математики имени акад. И.Н.Векуа
Тбилисского государственного университета, 1983, с. 260.

В книге рассматривается довольно широкий класс задач вычислительной математики и вычислительные процессы, которые приводят к построению решения данной задачи, анализируются источники погрешностей и оцениваются самые погрешности в предположении, что данная задача поставлена точно. Как кажется автору, во многих случаях осуществляется точно четыре источника погрешности: погрешность аппроксимации, происходящая от замены данной задачи другой, упрощенной задачей; погрешность искажения, связанная с тем, что данные упрощенной задачи вычисляются неточно; погрешность алгоритма, возникающая в тех случаях, когда алгоритм, примененный для решения упрощенной задачи, дает после конечного числа шагов лишь приближенное ее решение; погрешность округления, возникающая из-за того, что предписанные упомянутого алгоритма выполняются неточно. Для широкого круга вычислительных процессов даны оценки всех перечисленных погрешностей.

Книжка содержит 9 глав. В главе I дана постановка проблемы и введены два более простых примера. В главах II-VI исследуются погрешности линейных вычислительных процессов, в главах VII-IX тот же вопрос рассматривается для нелинейных процессов.

Книжка содержит расширенное изложение докладов, сделанных автором в Институте прикладной математики имени академика И.Н.Векуа в 1982 году.

Եւրոմոն շնորհակալ ժողովուրդ
գամտարանի շնորհակալ ժողովուրդներ
/հոյնը զննել/
տնօրէնի շնորհակալ ժողովուրդներ
տնօրէն, 1983

Редактор издательства Н. Девдариани
Подписано в печать 16.05.83
Листов 16,5
Учётно-изд. листов 18,31
Заказ 48 33 03803 1-й экз 290
Цена I руб. 30 коп.

36,56

Եւրոմոն շնորհակալ ժողովուրդի շնորհակալ ժողովուրդներ
330043, տնօրէն, 43.

Ротапринт: Института прикладной математики ТГУ
330043. Тбилиси, 43.

