

Л.Л. ГЛАЗЫРИНА, М.М. КАРЧЕВСКИЙ

# Введение в численные методы

Учебное пособие

Казань

2017

УДК 519.3  
ББК 22.311  
Г52

**Научный редактор**

доктор физико-математических наук М.Ф. Павлова

**Рецензенты:**

доктор физико-математических наук, профессор П.Г. Данилаев,  
доктор физико-математических наук, профессор В.С. Желтухин

**Глазырина Л.Л., Карчевский М.М.**

**Г 52 Введение в численные методы: учебное пособие.** — Казань: Казан. ун-т, 2017. — 122 с.

Излагаются основные принципы построения и исследования численных методов алгебры, анализа и решения дифференциальных уравнений. Приведены практические задания по методам, предложенным в пособии. Пособие предназначено для студентов, специализирующихся по математическим методам в экономике и информационным технологиям.

УДК 519.3  
ББК 22.311

© Глазырина Л.Л.  
Карчевский М.М., 2017

© Казанский университет,  
2017

---

---

## Оглавление

Предисловие . . . . .	4
<b>ГЛАВА 1. Численные методы алгебры . . . . .</b>	<b>5</b>
§ 1. Прямые методы решения систем линейных уравнений . . . . .	5
1. Метод Гаусса. . . . .	5
2. Метод отражений. . . . .	13
3. Метод Холесского. . . . .	15
4. Вычисление определителя и обратной матрицы. . . . .	16
5. Метод прогонки для систем с трехдиагональными матрицами. . . . .	16
§ 2. Итерационные методы решения систем линейных уравнений . . . . .	18
1. Методы Зейделя и Якоби. . . . .	18
2. Метод релаксации. . . . .	21
3. Пример решения задачи оптимизации итерационного параметра. . . . .	23
4. Итерационные методы вариационного типа. . . . .	25
§ 3. Методы решения алгебраической проблемы собственных значений . . . . .	27
1. Метод прямой итерации. . . . .	28
2. Метод обратной итерации. . . . .	29
3. Метод вращений (Якоби). . . . .	30
§ 4. Методы решения нелинейных уравнений . . . . .	33
1. Метод деления отрезка пополам. . . . .	33
2. Метод простой итерации. . . . .	34
3. Метод Ньютона. . . . .	37
4. Метод хорд. . . . .	37
5. Метод секущих. . . . .	39
§ 5. Методы решения систем нелинейных уравнений . . . . .	40
1. Метод простой итерации. . . . .	40
2. Метод Ньютона. . . . .	41
<b>ГЛАВА 2. Численные методы анализа . . . . .</b>	<b>42</b>
§ 1. Интерполирование функций . . . . .	42
1. Существование и единственность интерполяционного полинома. . . . .	42
2. Интерполяционный полином Лагранжа. . . . .	43
3. Интерполяционный полином Ньютона. . . . .	43
4. Оценка погрешности интерполирования. . . . .	45
5. Оптимальный выбор узлов интерполирования. Многочлены Чебышева. . . . .	47
6. Интерполирование с кратными узлами. . . . .	51
§ 2. Среднеквадратичное приближение функций . . . . .	53
1. Элемент наилучшего среднеквадратичного приближения. . . . .	53
2. Ортогональные полиномы. . . . .	56
§ 3. Приближенное вычисление интегралов . . . . .	60
1. Интерполяционные квадратурные формулы. . . . .	60

2.	Устойчивость квадратурных формул. . . . .	62
3.	Квадратурные формулы Ньютона — Котеса. . . . .	62
4.	Оценки точности простейших квадратурных формул Ньютона — Котеса. . . . .	64
5.	Составные квадратурные формулы. . . . .	67
6.	Квадратурные формулы типа Гаусса. . . . .	70
<b>ГЛАВА 3. Численные методы решения дифференциальных уравнений</b> . . . . . 75		
§ 1.	Численные методы решения задачи Коши . . . . .	75
1.	Метод, основанный на формуле Тейлора. . . . .	75
2.	Методы типа Рунге — Кутты. . . . .	77
3.	Элементы теории одношаговых методов решения задачи Коши. . . . .	79
4.	Методы типа Адамса. . . . .	83
§ 2.	Методы решения уравнений с частными производными . . . . .	86
1.	Сеточные методы решения краевых задач для обыкновенных дифференциальных уравнений. . . . .	86
2.	Вариационные методы. . . . .	90
3.	Метод Рунге. . . . .	92
4.	Метод конечных элементов. . . . .	94
5.	Метод конечных элементов для эллиптических уравнений. . . . .	96
6.	Итерационные методы решения сеточных уравнений. . . . .	102
7.	Разностные методы решения нестационарных задач математической физики. . . . .	103
<b>ГЛАВА 4. Практикум по численным методам</b> . . . . . 110		
§ 1.	Системы линейных уравнений . . . . .	110
§ 2.	Нелинейные уравнения . . . . .	111
§ 3.	Интерполирование функций . . . . .	114
§ 4.	Численное интегрирование . . . . .	116
§ 5.	Задача Коши для системы обыкновенных дифференциальных уравнений . . . . .	118
<b>Литература</b> . . . . .		121

---

---

## Предисловие

Пособие написано на основе лекций для студентов института вычислительной математики и информационных технологий КФУ, специализирующихся в области математических методов в экономике и информационных технологий. В пособии содержатся также типовые задания по основным разделам численных методов, которые могут быть использованы на практических занятиях. Предполагается, что читатель знаком со стандартными курсами математического анализа, линейной алгебры и аналитической геометрии, обыкновенных дифференциальных уравнений. Многие вопросы, затронутые в пособии, активно обсуждались с сотрудниками кафедры вычислительной математики Казанского федерального университета. Авторы выражают им свою искреннюю благодарность. Рукопись пособия была внимательно прочитана М.Ф.Павловой. Авторы постарались максимально учесть ее замечания. Авторы признательны Е.М.Федотову, оказавшему большую помощь при подготовке пособия к печати.

---

---

ГЛАВА 1  
Численные методы алгебры

§ 1. Прямые методы решения систем линейных уравнений

**1. Метод Гаусса.** Рассматривается система линейных алгебраических уравнений

$$Ax = b. \quad (1.1)$$

Здесь

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}, \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix}, \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \dots \\ b_n \end{pmatrix}.$$

Матрица  $A$  предполагается невырожденной, то есть  $\det A \neq 0$ . Поэтому система (1.1) однозначно разрешима при любой правой части, ее решение может быть выписано по формулам Крамера

$$x_i = \frac{\Delta_i}{\Delta}, \quad i = 1, 2, \dots, n,$$

где  $\Delta = \det A$  — определитель матрицы  $A$ <sup>1)</sup>, а  $\Delta_i$  — это определитель, получающийся из определителя матрицы  $A$  заменой  $i$ -го столбца столбцом свободных членов системы (1.1).

Казалось бы, формулы Крамера полностью решают задачу построения решения системы линейных уравнений, однако на практике они не используются. Это объясняется следующим. Напомним, что

$$\det A = \sum_{(\alpha_1, \alpha_2, \dots, \alpha_n)} \pm a_{1\alpha_1} a_{2\alpha_2} \dots a_{n\alpha_n}, \quad (1.2)$$

где  $(\alpha_1, \alpha_2, \dots, \alpha_n)$  — перестановка символов  $1, 2, \dots, n$ . Число слагаемых в сумме равно  $n!$ , поэтому непосредственное вычисление определителя требует  $nn!$  арифметических операций, что уже при  $n = 30$  недоступно даже для самых мощных ЭВМ. Кроме того, из-за большого числа сомножителей в слагаемых (1.2) непосредственные вычисления по формуле (1.2) могут приводить к переполнению разрядной

---

<sup>1)</sup>Определитель матрицы  $A$  обозначается также через  $|A|$ .

сетки ЭВМ или к сильному накоплению погрешностей округления. Поэтому для решения систем уравнений применяют другие, более экономичные и устойчивые по отношению к погрешностям округления методы.

**1.1.** Начнем с самого простого и наиболее распространенного метода — метода Гаусса. Запишем систему (1.1) в индексной форме:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= b_1, \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= b_2, \\ \dots & \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n &= b_n. \end{aligned} \quad (1.3)$$

Предположим, что  $a_{11} \neq 0$  и поделим на это число первое уравнение системы. Получим:

$$x_1 + a_{12}^{(1)}x_2 + \dots + a_{1n}^{(1)}x_n = b_1^{(1)}.$$

Умножим теперь первое уравнение системы на  $a_{21}$  и вычтем из второго. Аналогично преобразуем остальные уравнения системы. В результате получим систему уравнений, эквивалентную исходной:

$$\begin{aligned} x_1 + a_{12}^{(1)}x_2 + \dots + a_{1n}^{(1)}x_n &= b_1^{(1)}, \\ a_{22}^{(1)}x_2 + \dots + a_{2n}^{(1)}x_n &= b_2^{(1)}, \\ \dots & \\ a_{n2}^{(1)}x_2 + \dots + a_{nn}^{(1)}x_n &= b_n^{(1)}, \end{aligned} \quad (1.4)$$

где

$$\begin{aligned} a_{1j}^{(1)} &= a_{1j}/a_{11}, \quad j = 2, \dots, n, \quad b_1^{(1)} = b_1/a_{11}, \quad a_{ij}^{(1)} = a_{ij} - a_{1j}^{(1)}a_{i1}, \quad (1.5) \\ i &= 2, \dots, n, \quad j = 2, \dots, n, \quad b_i^{(1)} = b_i - b_1a_{i1}, \quad i = 2, \dots, n. \end{aligned}$$

Отбросим первое уравнение системы (1.4), а для оставшихся — выполним преобразования, такие же, как для системы (1.3). Получим:

$$\begin{aligned} x_1 + a_{12}^{(1)}x_2 + a_{13}^{(1)}x_3 + \dots + a_{1n}^{(1)}x_n &= b_1^{(1)}, \\ x_2 + a_{23}^{(2)}x_3 + \dots + a_{2n}^{(2)}x_n &= b_2^{(2)}, \\ a_{33}^{(2)}x_3 + \dots + a_{3n}^{(2)}x_n &= b_3^{(2)}, \\ \dots & \\ a_{n3}^{(2)}x_3 + \dots + a_{nn}^{(2)}x_n &= b_n^{(2)}. \end{aligned} \quad (1.6)$$





**1.3.** Условия применимости метода Гаусса. Описанный метод решения системы линейных уравнений может быть реализован лишь в том случае, когда все числа  $a_{11}, a_{22}^{(1)}, \dots, a_{nn}^{(n-1)}$ , называемые ведущими или главными элементами метода Гаусса, отличны от нуля. Выделим класс матриц, для которых это условие выполняется. Пусть

$$A_1 = a_{11}, \quad A_2 = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix}, \quad \dots, \quad A_n = \begin{vmatrix} a_{11} & a_{22} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{vmatrix}$$

есть главные миноры матрицы  $A$ .

**Теорема 1.1.** *Для того, чтобы все ведущие элементы метода Гаусса были отличны от нуля, необходимо и достаточно, чтобы все главные миноры матрицы  $A$  были ненулевыми.*

**ДОКАЗАТЕЛЬСТВО.** Пусть все главные миноры матрицы  $A$  отличны от нуля. Покажем, что тогда все ведущие элементы метода Гаусса не равны нулю. Имеем, в частности,  $a_{11} = A_1 \neq 0$ . Применяя преобразования, выполнявшиеся при проведении прямого хода метода Гаусса (см. формулы (1.5)), получим

$$A_2 = a_{11} \begin{vmatrix} 1 & \frac{a_{12}}{a_{11}} \\ 0 & a_{22} - \frac{a_{12}a_{21}}{a_{11}} \end{vmatrix} = a_{11}a_{22}^{(1)},$$

следовательно,  $a_{22}^{(1)} \neq 0$ . Пусть уже доказано, что  $a_{11}, a_{22}^{(1)}, \dots, a_{k-1,k-1}^{(k-2)}$  не равны нулю. Тогда, приводя минор  $A_k$  к треугольному виду при помощи преобразований прямого хода метода Гаусса, получим

$$A_k = a_{11}a_{22}^{(1)} \dots a_{k-1,k-1}^{(k-2)} \begin{vmatrix} 1 & a_{12}^{(1)} & \dots & a_{1k}^{(1)} \\ 0 & 1 & \dots & a_{2k}^{(2)} \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & a_{kk}^{(k-1)} \end{vmatrix} =$$

$$= a_{11}a_{22}^{(1)} \dots a_{k-1,k-1}^{(k-2)} a_{kk}^{(k-1)}, \quad (1.9)$$

следовательно,  $a_{kk}^{(k-1)} \neq 0$ . Обратное утверждение теоремы есть очевидное следствие соотношения (1.9).  $\square^1$

<sup>1)</sup>Значком  $\square$ , как обычно, отмечаем конец доказательства.



$$\leq |x_i| \sum_{j=1, j \neq i}^k |a_{ij}| \leq |x_i| \sum_{j=1, j \neq i}^n |a_{ij}|,$$

поэтому  $|a_{ii}| \leq \sum_{j=1, j \neq i}^n |a_{ij}|$ , что противоречит условию (1.10).  $\square$

**1.4.** Метод Гаусса с выбором главных элементов. В тех случаях, когда заранее трудно установить, что все ведущие элементы метода Гаусса отличны от нуля, применяются модификации метода Гаусса, позволяющие избежать деления на нуль в случае произвольной невырожденной матрицы  $A$ . Опишем вариант такого метода, называемый методом Гаусса с выбором главного элемента по строке.

Найдем максимальный по модулю элемент первой строки матрицы  $A$ . Пусть это есть  $a_{1j}$ . Ясно, что  $a_{1j} \neq 0$ , иначе все элементы первой строки матрицы равны нулю и матрица оказывается вопреки нашему предположению вырожденной. Поменяем теперь местами первый и  $j$ -й столбцы системы (1.3) и выполним первый шаг прямого хода метода Гаусса. При этом будет выполняться деление на  $a_{1j} \neq 0$ . Затем проводится поиск максимального по модулю элемента во второй строке системы вида (1.4). Он вновь оказывается ненулевым, поскольку матрица  $A$  невырождена. Выполняется перестановка соответствующего столбца системы (1.4) с ее вторым столбцом, затем — второй шаг прямого хода метода Гаусса. Аналогичные вычисления проводятся до приведения системы (1.3) к треугольному виду.

При программной реализации этого метода перестановки столбцов заменяются соответствующими перенумерациями неизвестных.

Ясно, что аналогичный алгоритм можно построить, проводя каждый раз поиск максимального элемента в соответствующем столбце матрицы. Можно было бы проводить поиск максимального элемента и по всей матрице.

Отметим, что выбор главного элемента не только обеспечивает реализуемость метода Гаусса при любой невырожденной матрице  $A$ , но и уменьшает влияние погрешностей округления, так как приводит к делению на большие числа, чем в основном алгоритме метода Гаусса. Понятно, что это требует дополнительных затрат на поиск максимального элемента и запоминание перенумераций неизвестных. Особенно значительными эти затраты становятся при поиске ведущего элемента по всей матрице. Однако этот вариант метода Гаусса оказывается наиболее устойчивым по отношению к погрешностям округления.

**1.5.** Метод Гаусса и разложение матрицы на треугольные множители. На практике чаще всего метод Гаусса реализуется в форме разложения матрицы на треугольные множители. Возможность такого разложения гарантирует

**Теорема 1.3.** Пусть все главные миноры матрицы  $A$  отличны от нуля. Тогда существуют невырожденные нижняя и верхняя треугольные матрицы

$$L = \begin{pmatrix} l_{11} & 0 & \dots & 0 \\ l_{12} & l_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{pmatrix}, \quad U = \begin{pmatrix} 1 & u_{12} & \dots & u_{1n} \\ 0 & 1 & \dots & u_{n2} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

(все диагональные элементы матрицы  $U$  равны единице) такие, что  $A = LU$ .

Прежде чем доказывать теорему, заметим, что если матрицы  $L$ ,  $U$  построены, то решение системы (1.1) может быть выполнено следующим образом. Запишем систему (1.1) в виде

$$LUx = b.$$

Положим

$$Ux = y. \quad (1.12)$$

Тогда

$$Ly = b. \quad (1.13)$$

Система (1.13) — система с треугольной матрицей. Ее решение может быть построено при помощи обратного хода метода Гаусса. Система (1.12) при уже найденном векторе  $y$  также может быть решена при помощи обратного хода метода Гаусса (заметим только, что неизвестные  $x_1, \dots, x_n$  находятся при этом начиная с первого). Таким образом, если матрицы  $L, U$  построены, то для решения системы (1.1) придется потратить примерно  $2n^2$  арифметических операций.

Особенно полезна такая модификация метода Гаусса при необходимости решать множество систем с одной и той же матрицей и различными правыми частями.

На практике такая ситуация возникает довольно часто. Объясняется это тем, что матрица  $A$ , обычно, описывает структуру некоторой системы (механической, экономической и т. д.), а вектор  $b$  характеризует внешние воздействия на эту систему. Поэтому при изучении реакции системы на различные внешние воздействия приходится решать множество систем с одной и той же матрицей.

**ДОКАЗАТЕЛЬСТВО** теоремы 1.3. Используем индукцию по порядку матрицы  $A$ . Для матрицы порядка  $k = 1$  утверждение теоремы тривиально. Пусть оно верно для некоторого  $k \geq 1$ . Докажем, что тогда оно будет верным и для матрицы порядка  $k + 1$ . Запишем матрицу  $A_{k+1}$  порядка  $k + 1$  в блочном виде

$$A_{k+1} = \begin{pmatrix} A_k & a_k \\ b_k & a_{k+1,k+1} \end{pmatrix}. \quad (1.14)$$

Здесь  $A_k$  — матрица порядка  $k$ ,  $a_k$  — столбец длины  $k$ ,  $b_k$  — строка длины  $k$ . Матрица  $A_k$  по условию теоремы невырождена и по индуктивному предположению может быть представлена в виде  $A_k = L_k U_k$ , где  $L_k$  — нижняя треугольная матрица с ненулевыми элементами на диагонали, а  $U_k$  — верхняя треугольная матрица с единичными диагональными элементами. Будем искать матрицу  $A_{k+1}$  в виде произведения двух блочных матриц:

$$A_{k+1} = \begin{pmatrix} L_k & 0 \\ l_k & l_{k+1,k+1} \end{pmatrix} \begin{pmatrix} U_k & u_k \\ 0 & 1 \end{pmatrix}. \quad (1.15)$$

Здесь  $l_k$  — искомая строка длины  $k$ ,  $u_k$  — искомый столбец длины  $k$ . Подсчитывая произведение сомножителей в правой части равенства (1.15) и приравнивая результат поблочно матрице  $A_{k+1}$  (см. (1.14)), получим

$$L_k U_k = A_k, \quad l_k U_k = b_k, \quad L_k u_k = a_k, \quad l_k u_k + l_{k+1,k+1} = a_{k+1,k+1}.$$

Отсюда можно найти  $l_k$ ,  $u_k$ , решая системы уравнений

$$U_k^T l_k^T = b_k^T, \quad (1.16)$$

$$L_k u_k = a_k, \quad (1.17)$$

и затем вычислить

$$l_{k+1,k+1} = a_{k+1,k+1} - l_k u_k. \quad (1.18)$$

При этом  $l_{k+1,k+1}$  не может обратиться в нуль, так как вследствие равенства (1.15) имеем  $|A_{k+1}| = |L_k| l_{k+1,k+1}$ , а определитель  $|A_{k+1}|$  по сделанному нами предположению отличен от нуля. Таким образом, искомое треугольное разложение матрицы  $A_{k+1}$  построено.  $\square$

Доказательство теоремы 1.3, фактически, дает алгоритм построения матриц  $L$  и  $U$ : нужно последовательно строить разложения диагональных блоков матрицы порядков  $k = 1, 2, \dots, n$ , используя формулы (1.16)–(1.18);  $k$ -й шаг такого алгоритма состоит в решении систем с треугольными матрицами и требует примерно  $2k^2$  операций.

Таким образом, построение матриц  $L, U$  требует примерно  $2n^3/3$  операций.

Нетрудно проверить, что матрица  $U$  совпадает с матрицей системы (1.7).

Если нет необходимости сохранять матрицу  $A$ , то программу вычислений можно организовать так, что элементы матриц  $L, U$  будут последовательно замещать соответствующие элементы матрицы  $A$ .

Аналогично методу Гаусса с выбором главных элементов можно строить алгоритмы разложения на треугольные множители, применимые для любой невырожденной матрицы.

**2. Метод отражений.** Вместо разложения матрицы на треугольные множители при решении системы уравнений можно использовать представление матрицы  $A$  в виде  $A = QU$ , где  $Q$  — ортогональная матрица, т. е. матрица, удовлетворяющая условию  $QQ^T = E$ , где  $E$  — единичная матрица, а  $U$  — верхняя треугольная матрица. Если такое разложение получено, то решение системы (1.1) сводится к последовательному решению систем  $Qy = b$ ,  $Ux = y$ . При этом ясно, что  $y = Q^T b$ , а  $x$  находится при помощи обратного хода метода Гаусса. Если матрица  $Q$  известна, то вычисление вектора  $y$  требует  $2n^2$  арифметических операций. Таким образом, после получения разложения матрицы  $A$  решения системы уравнений требует примерно  $3n^2$  операций.

Для построения указанного разложения будем использовать матрицы отражения, т. е. матрицы вида  $R = E - 2ww^T$ , где  $w$  — единичный вектор:  $(w, w) = |w|^2 = 1$ . Матрица  $R$  при любом  $|w| = 1$  симметрична и ортогональна. Действительно,

$$R^T = R, \quad R^2 = E - 4ww^T + 4ww^Tww^T = E,$$

так как  $w^T w = |w|^2 = 1$ . Заметим, далее, что

$$Rw = w - 2ww^T w = -w, \quad Rz = z - 2ww^T z = z, \quad (2.1)$$

если  $w^T z = (w, z) = 0$ , т. е. векторы  $w$  и  $z$  ортогональны.

Пусть теперь  $x$  — произвольный вектор. По теореме об ортогональном разложении евклидова пространства он однозначно представим в виде  $x = \alpha w + z$ , где  $\alpha$  некоторое число,  $z$  — некоторый вектор, ортогональный  $w$ . Из равенств (2.1) вытекает, что  $Rx = -\alpha w + z$ .

Можно сказать, таким образом, что матрица  $R$  выполняет отражение вектора  $x$  относительно  $(n - 1)$ -мерной гиперплоскости, ортогональной вектору  $w$ . Это свойство матрицы  $R$  и позволяет называть ее матрицей отражения.

Рассмотрим теперь следующую задачу. Даны ненулевой вектор  $x$  и единичный вектор  $e$ . Требуется построить матрицу отражения  $R$ , такую, что  $Rx = \mu e$ , где  $\mu$  — число (ясно, что  $|\mu| = |x|$ , поскольку матрица  $R$  ортогональна).

Нетрудно видеть (сделайте чертеж!), что решение задачи — матрица отражения, определяемая вектором  $w = (x - |x|e)/|x - |x|e|$  или вектором  $w = (x + |x|e)/|x + |x|e|$ .

При вычислениях для минимизации погрешностей округления следует выбрать вектор  $w$  с бóльшим знаменателем.

**Теорема 2.1.** Пусть  $A$  — произвольная невырожденная матрица. Тогда существует ортогональная матрица  $Q$  такая, что  $A = QU$ , где  $U$  — верхняя треугольная матрица с ненулевыми диагональными элементами.

**ДОКАЗАТЕЛЬСТВО.** Пусть  $a^1$  — первый столбец матрицы  $A$ ,  $e^1 = (1, 0, \dots, 0)^T$ ,  $w^1 = (a^1 \pm |a^1|e^1)/|a^1 \pm |a^1|e^1|$  (знак выбирается, как указано выше),  $R_1 = E - 2w^1w^{1T}$ . образуем матрицу  $A_1 = (a_{ij}^{(1)})_{i,j=1}^n = R_1A$ . Ясно, что первый столбец этой матрицы коллинеарен вектору  $e^1$  и, следовательно, имеет вид  $(a_{11}^{(1)}, 0, \dots, 0)^T$ , причем  $|a_{11}^{(1)}| = |a^1| \neq 0$ .

Рассмотрим столбец  $\tilde{a}^2 = (a_{22}^{(1)}, a_{32}^{(1)}, \dots, a_{n2}^{(1)})^T$  длины  $n - 1$ . Понятно, что  $\tilde{a}^2 \neq 0$ , так как в противном случае  $|A_1| = 0$ , но, с другой стороны, поскольку определитель ортогональной матрицы равен  $\pm 1$ , то  $|A_1| = \pm|A|$ , а матрица  $A$  по предположению невырождена. Пусть  $\tilde{w} = (\tilde{a}^2 \pm |\tilde{a}^2|e^1)/|\tilde{a}^2 \pm |\tilde{a}^2|e^1|$ , где  $e^1 = (1, 0, \dots, 0)$  — вектор длины  $n - 1$ . Положим  $\tilde{R}_2 = \tilde{E} - 2\tilde{w}\tilde{w}^T$ , где  $\tilde{E}$  — единичная матрица размерности  $n - 1$ , и образуем ортогональную матрицу

$$R_2 = \begin{pmatrix} 1 & 0 \\ 0^T & \tilde{R}_2 \end{pmatrix}.$$

Здесь  $0$  — нулевая строка длины  $n - 1$ . Пусть  $A_2 = (a_{ij}^{(2)})_{i,j=1}^n = R_2A_1$ . Нетрудно видеть, что первые столбцы матриц  $A_1, A_2$  совпадают, а второй столбец матрицы  $A_2$  имеет вид  $(a_{12}^{(2)}, a_{22}^{(2)}, 0, \dots, 0)^T$ , причем  $|a_{22}^{(2)}| = |\tilde{a}^2| \neq 0$ .

Выполняя аналогичные построения, получим ортогональные матрицы  $R_3, R_4, \dots, R_n$  такие, что  $R_n R_{n-1} \dots R_1 A = A_n$ , где  $A_n$  — верхняя треугольная матрица с ненулевыми элементами  $a_{11}^{(1)}, a_{22}^{(2)}, \dots, a_{nn}^{(n)}$  на диагонали, следовательно,  $A = QU$ , где  $U = A_n, Q = R_1^T R_2^T \dots R_n^T$  суть ортогональная матрица.  $\square$

Доказательство теоремы 2.1, фактически, дает алгоритм построения матриц  $Q, U$ . Нетрудно убедиться, что их вычисление требует приблизительно в два раза больше арифметических операций, чем разложение матрицы на треугольные множители. Важным положительным качеством описанного алгоритма является возможность его применения для произвольной невырожденной матрицы без какой-либо перенумерации строк или столбцов.

**3. Метод Холецкого.** В случае, когда матрица системы линейных уравнений симметрична и положительно определена, можно добиться существенного сокращения числа операций и памяти, необходимых для вычисления решения. В основе соответствующего алгоритма лежит

**Теорема 3.1.** Пусть матрица  $A$  симметрична и положительно определена. Тогда существует нижняя треугольная матрица  $L$  с положительными элементами на диагонали такая, что  $A = LL^T$ .

ДОКАЗАТЕЛЬСТВО. Используем индукцию по порядку матрицы  $A$ . Для матрицы порядка  $k = 1$  имеем тривиальное равенство  $A_1 = a_{11} = \sqrt{a_{11}}\sqrt{a_{11}}$ . Пусть нужное разложение получено для некоторого  $k \geq 1$ . Покажем, как его построить для матрицы порядка  $k+1$ . Запишем матрицу  $A_{k+1}$  как блочную:

$$A_{k+1} = \begin{pmatrix} A_k & a_k \\ a_k^T & a_{k+1,k+1} \end{pmatrix}.$$

В силу предположения индукции  $A_k = L_k L_k^T$ , где  $L_k$  — нижняя треугольная матрица с положительными элементами на диагонали. Будем искать разложение матрицы  $A$  на треугольные множители в виде

$$A_{k+1} = L_{k+1} L_{k+1}^T = \begin{pmatrix} L_k & 0 \\ l_k^T & l_{k+1,k+1} \end{pmatrix} \begin{pmatrix} L_k^T & l_k \\ 0 & l_{k+1,k+1} \end{pmatrix}. \quad (3.1)$$

Выполняя умножение в правой части последнего равенства и сравнивая поблочно результат с матрицей  $A_{k+1}$ , получим систему линейных уравнений

$$L_k l_k = a_k \quad (3.2)$$

для определения вектора  $l_k$  и уравнение  $l_k^T l_k + l_{k+1,k+1}^2 = a_{k+1,k+1}$  для элемента  $l_{k+1,k+1}$ . Можно считать, что  $l_{k+1,k+1} > 0$ , так как вследствие (3.1) имеем:  $|A_{k+1}| = |A_k| l_{k+1,k+1}^2$ , причем  $|A_k|, |A_{k+1}| > 0$ , так как по условию матрицы  $A_k, A_{k+1}$  положительно определены. Таким образом, для построения матрицы  $L_{k+1}$  нужно решить систему уравнений (3.2) с треугольной матрицей, а затем вычислить  $l_{k+1,k+1}$  по формуле  $l_{k+1,k+1} = \sqrt{a_{k+1,k+1} - l_k^T l_k}$ .  $\square$



Доказательство теоремы 3.1, фактически, описывает алгоритм разложения на треугольные множители произвольной симметричной положительно определенной матрицы. Нетрудно видеть, что его реализация по затратам памяти и объему вычислений оказывается в два раза более экономичной, чем разложение на треугольные множители произвольной невырожденной матрицы.

После того, как матрица  $L$  построена, решение задачи (1.1) сводится к последовательному решению систем уравнений  $Ly = b$ ,  $L^T x = y$  с треугольными матрицами.

**4. Вычисление определителя и обратной матрицы.** Изложенные выше методы решения систем линейных уравнений позволяют решать и эти задачи.

Так, если применяется метод Гаусса, то попутно можно вычислить и  $|A| = a_{11}a_{22}^{(1)} \dots a_{nn}^{(n-1)}$  (см. (1.9)). Аналогичные формулы нетрудно написать и в тех случаях, когда строится разложение матрицы  $A$  на множители. Надо, однако, иметь в виду, что непосредственные вычисления по этим формулам, как правило, оказываются невозможными: из-за большого числа сомножителей определитель (или результат промежуточных вычислений) зачастую либо слишком велик, либо, наоборот, слишком мал. Приходится писать специальные программы, позволяющие отдельно подсчитывать мантиссу и порядок определителя.

Построение обратной матрицы сводится к решению  $n$  систем линейных уравнений с одной и той же матрицей  $A$  и различными правыми частями. Действительно, обозначим матрицу  $A^{-1}$  через  $X$ . Тогда  $AX = E$ . Осталось записать это равенство подробнее:

$$Ax^k = e^k, \quad k = 1, 2, \dots, n.$$

Здесь  $x^k$  —  $k$ -й столбец матрицы  $X$ ,

$$e^k = (\underbrace{0, \dots, 0}_{k-1}, 1, \underbrace{0, \dots, 0}_{n-k}).$$

**5. Метод прогонки для систем с трехдиагональными матрицами.** В приложениях довольно часто возникают системы уравнений с матрицами, большинство элементов которых — нули. Это так называемые разреженные матрицы. Процесс исключения неизвестных в таких системах (или разложение матриц на треугольные множители) во многих практически важных ситуациях удается организовать так, чтобы существенно сократить память и объем вычислений.



Описанный алгоритм носит название метода прогонки. Понятно, что это — метод Гаусса, записанный применительно к случаю трехдиагональной системы уравнений, причем процесс вычислений  $P_i, Q_i$  (прямой ход метода прогонки) соответствует прямому ходу метода Гаусса, а вычисления по формулам (5.4) (обратный ход метода прогонки) соответствуют обратному ходу метода Гаусса.

Нетрудно подсчитать необходимые затраты: требуется примерно  $8n$  арифметических операций и не более  $6n$  ячеек памяти.

Метод может быть реализован, когда все знаменатели в формулах (5.3), (5.5) отличны от нуля. Учитывая связь метода прогонки с методом Гаусса, можно сказать, что данное условие выполнено, например, когда матрица системы (5.1) — матрица с диагональным преобладанием, т. е.  $|c_1| < |b_1|$ ,  $|a_n| < |b_n|$ ,  $|a_i| + |c_i| < |b_i|$ ,  $i = 2, \dots, n - 1$ .

## § 2. Итерационные методы решения систем линейных уравнений

Рассмотренные выше методы решения систем линейных алгебраических уравнений принято называть прямыми методами. Все они характеризуются тем, что если пренебречь ошибками округления, то решение системы может быть получено за конечное число арифметических операций (зависящее лишь от порядка системы).

При реализации прямых методов важно, чтобы все данные располагались в оперативной (быстрой) памяти компьютера. Если порядок системы настолько велик, что ее матрица может быть размещена только во внешней (медленной) памяти, например, на жестком диске, то время, затрачиваемое на решение системы, существенно увеличивается.

Поэтому для больших систем предпочтительнее оказываются итерационные методы. Основная идея этих методов состоит в построении последовательности векторов  $x^k$ ,  $k = 1, 2, \dots$ , сходящейся к решению системы (1.1). За приближенное решение принимается вектор  $x^k$  при достаточно большом  $k$ . При реализации итерационных методов, обычно, достаточно уметь вычислять вектор  $Ax$  при любом заданном векторе  $x$ .

**1. Методы Зейделя и Якоби.** Будем считать, что все диагональные элементы матрицы  $A$  отличны от нуля, и перепишем систему (1.1), разрешая каждое уравнение относительно переменной,

стоящей на диагонали:

$$x_i = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j + \frac{b_i}{a_{ii}}, \quad i = 1, 2, \dots, n. \quad (1.1)$$

Выберем некоторое начальное приближение  $x^0 = (x_1^0, x_2^0, \dots, x_n^0)^T$  и построим последовательность векторов  $x^1, x^2, \dots$ , определяя вектор  $x^{k+1}$  по уже найденному вектору  $x^k$  при помощи соотношений:

$$x_i^{k+1} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j^k - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j^k + \frac{b_i}{a_{ii}}, \quad i = 1, 2, \dots, n. \quad (1.2)$$

Формулы (1.2) определяют итерационный метод решения системы (1.1), называемый методом Якоби или методом простой итерации.

Укажем легко проверяемое достаточное условие сходимости этого метода.

**Теорема 1.1.** Пусть матрица  $A$  — матрица с диагональным преобладанием. Тогда итерационный метод Якоби сходится при любом начальном приближении  $x^0$ .

**ДОКАЗАТЕЛЬСТВО.** Заметим, прежде всего, что условие диагонального преобладания (1.10) означает, что

$$\max_{1 \leq i \leq n} \sum_{j=1, j \neq i}^n \frac{|a_{ij}|}{|a_{ii}|} = q < 1. \quad (1.3)$$

Пусть  $x$  — решение системы уравнений (1.1). Здесь и всюду в дальнейшем погрешность метода на  $k$ -м шаге итераций, т. е. вектор  $x^k - x$ , будем обозначать через  $z^k$ . Вычитая почленно из равенства (1.2) равенство (1.1), получим

$$z_i^{k+1} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} z_j^k - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} z_j^k, \quad i = 1, 2, \dots, n,$$

следовательно,

$$\begin{aligned} |z_i^{k+1}| &\leq \sum_{j=1}^{i-1} \frac{|a_{ij}|}{|a_{ii}|} |z_j^k| + \sum_{j=i+1}^n \frac{|a_{ij}|}{|a_{ii}|} |z_j^k| \leq \left( \sum_{j=1}^{i-1} \frac{|a_{ij}|}{|a_{ii}|} + \sum_{j=i+1}^n \frac{|a_{ij}|}{|a_{ii}|} \right) \max_{1 \leq j \leq n} |z_j^k| = \\ &= q \max_{1 \leq j \leq n} |z_j^k|, \end{aligned}$$

$i = 1, 2, \dots, n$ , откуда вытекает, что

$$\max_{1 \leq j \leq n} |z_j^{k+1}| \leq q \max_{1 \leq j \leq n} |z_j^k|$$

для любого  $k = 0, 1, \dots$ , поэтому

$$\max_{1 \leq j \leq n} |z_j^k| \leq q^k \max_{1 \leq j \leq n} |z_j^0| \rightarrow 0 \quad (1.4)$$

при  $k \rightarrow \infty$ , поскольку  $0 < q < 1$ , а это и означает, что  $x^k \rightarrow x$ .  $\square$

Оценка (1.4) показывает, что, чем меньше  $q$ , тем быстрее сходится метод простой итерации.

Формулы (1.2) допускают естественную модификацию. Именно, при вычислении  $x_i^{k+1}$  будем использовать уже найденные компоненты вектора  $x^{k+1}$ , то есть  $x_1^{k+1}, x_2^{k+1}, \dots, x_{i-1}^{k+1}$ . В результате приходим к итерационному методу Зейделя:

$$x_i^{k+1} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j^{k+1} - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j^k + \frac{b_i}{a_{ii}}, \quad i = 1, 2, \dots, n, \quad (1.5)$$

$k = 0, 1, \dots$

Метод Зейделя позволяет более экономно расходовать память, поскольку в данном случае вновь получаемые компоненты вектора  $x^{k+1}$  можно размещать на месте соответствующих компонент вектора  $x^k$ , в то время как при реализации метода Якоби все компоненты векторов  $x^k, x^{k+1}$  должны одновременно находиться в памяти компьютера.

Достаточное условие сходимости и оценку скорости сходимости метода Зейделя дает

**Теорема 1.2.** Пусть матрица  $A$  — матрица с диагональным преобладанием. Тогда метод Зейделя сходится при любом начальном приближении  $x^0$ ; справедлива оценка скорости сходимости:

$$\max_{1 \leq j \leq n} |z_j^k| \leq q^k \max_{1 \leq j \leq n} |z_j^0|. \quad (1.6)$$

**ДОКАЗАТЕЛЬСТВО.** Вычитая почленно из равенства (6.2) равенство (1.1), получим

$$z_i^{k+1} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} z_j^{k+1} - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} z_j^k, \quad i = 1, 2, \dots, n. \quad (1.7)$$

Пусть  $\max_{1 \leq j \leq n} |z_j^{k+1}| = |z_l^{k+1}|$ . Из  $l$ -го уравнения системы (1.7) вытекает, что

$$|z_l^{k+1}| \leq \alpha_l \max_{1 \leq j \leq n} |z_j^{k+1}| + \beta_l \max_{1 \leq j \leq n} |z_j^k|,$$

где

$$\alpha_l = \sum_{j=1}^{l-1} \frac{|a_{lj}|}{|a_{ll}|}, \quad \beta_l = \sum_{j=l+1}^n \frac{|a_{lj}|}{|a_{ll}|},$$

следовательно,

$$\max_{1 \leq j \leq n} |z_j^{k+1}| \leq \frac{\beta_l}{1 - \alpha_l} \max_{1 \leq j \leq n} |z_j^k|.$$

Из условия (1.3) получаем, что  $\alpha_l + \beta_l \leq q < 1$ , но тогда и  $q\alpha_l + \beta_l \leq q$ , таким образом,  $\beta_l / (1 - \alpha_l) \leq q$ .  $\square$

**2. Метод релаксации.** Во многих ситуациях существенного ускорения сходимости можно добиться за счет введения так называемого итерационного параметра. Рассмотрим итерационный процесс

$$x_i^{k+1} = (1 - \omega)x_i^k + \omega \left( - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j^{k+1} - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j^k + \frac{b_i}{a_{ii}} \right), \quad (2.1)$$

$i = 1, 2, \dots, n, k = 0, 1, \dots$  Этот метод называется методом релаксации, число  $\omega$  — релаксационным параметром. При  $\omega = 1$  метод переходит в метод Зейделя.

Ясно, что по затратам памяти и объему вычислений на каждом шаге итераций метод релаксации не отличается от метода Зейделя.

Мы исследуем сходимость метода релаксации в случае, когда матрица  $A$  симметрична и положительно определена. С этой целью перепишем его в матричном виде. Обозначим через  $L$  нижнюю треугольную матрицу с нулевой главной диагональю; элементы, стоящие под главной диагональю матрицы  $L$ , совпадают с соответствующими элементами матрицы  $A$ . Через  $D$  обозначим диагональную матрицу, на диагонали которой стоят диагональные элементы матрицы  $A$ . Понятно, что  $A = L + D + L^T$ . Нетрудно убедиться, что равенства (2.1) с учетом введенных обозначений принимают вид:

$$Dx^{k+1} = (1 - \omega)Dx^k + \omega(-Lx^{k+1} - L^T x^k + b).$$

После элементарных преобразований получим, что

$$B \frac{x^{k+1} - x^k}{\omega} + Ax^k = b, \quad (2.2)$$

где  $B = D + \omega L$ .

Нам потребуется следующая общая теорема, полезная при исследовании многих итерационных методов.

**Теорема 2.1.** Пусть матрица  $A$  симметрична и положительно определена. Тогда итерационный метод

$$B \frac{x^{k+1} - x^k}{\tau} + Ax^k = b, k = 0, 1, \dots \quad (2.3)$$

где  $\tau > 0$ , сходится при любом начальном приближении  $x^0$ , если

$$(Bx, x) > \frac{\tau}{2}(Ax, x) \quad \forall x \neq 0. \quad (2.4)$$

ДОКАЗАТЕЛЬСТВО. Если  $x$  — решение уравнения (1.1), то, очевидно,

$$B \frac{x - x}{\tau} + Ax = b. \quad (2.5)$$

Вычитая почленно из равенства (2.5) равенство (2.3), получим

$$B \frac{z^{k+1} - z^k}{\tau} + Az^k = 0. \quad (2.6)$$

Используя тривиальное равенство  $z^k = \frac{z^{k+1} + z^k}{2} - \frac{z^{k+1} - z^k}{2}$ , преобразуем (2.6) к виду

$$\left( B - \frac{\tau}{2}A \right) \frac{z^{k+1} - z^k}{\tau} + \frac{1}{2}A(z^{k+1} + z^k) = 0. \quad (2.7)$$

Умножим обе части равенства (2.7) скалярно на вектор  $2(z^{k+1} - z^k)$ . После элементарных преобразований с учетом симметрии матрицы  $A$  получим:

$$2\tau \left( \left( B - \frac{\tau}{2}A \right) \frac{z^{k+1} - z^k}{\tau}, \frac{z^{k+1} - z^k}{\tau} \right) + (Az^{k+1}, z^{k+1}) - (Az^k, z^k) = 0. \quad (2.8)$$

Вследствие условия (2.4) первое слагаемое в правой части (2.8) неотрицательно, поэтому  $(Az^{k+1}, z^{k+1}) \leq (Az^k, z^k)$ , т. е. числовая последовательность  $(Az^k, z^k)$  не возрастает. Кроме того, она ограничена снизу нулем, так как матрица  $A$  положительно определена. Таким образом, последовательность  $(Az^k, z^k)$  имеет предел. Отсюда вытекает, что  $(Az^{k+1}, z^{k+1}) - (Az^k, z^k) \rightarrow 0$  при  $k \rightarrow \infty$ , следовательно, и

$$\left( \left( B - \frac{\tau}{2}A \right) (z^{k+1} - z^k), z^{k+1} - z^k \right) \rightarrow 0 \text{ при } k \rightarrow \infty,$$

а поскольку матрица  $B - (\tau/2)A$  положительно определена (см. условие (2.4)), то  $z^{k+1} - z^k \rightarrow 0$  при  $k \rightarrow \infty$ . Используя теперь уравнение (2.6) и невырожденность матрицы  $A$ , получим, что  $z^k \rightarrow 0$  при  $k \rightarrow \infty$ .  $\square$

С использованием теоремы 2.1 просто доказывается

**Теорема 2.2.** Пусть матрица  $A$  положительно определена, параметр релаксации удовлетворяет условию:  $0 < \omega < 2$ . Тогда метод релаксации сходится при любом начальном приближении  $x^0$ .

ДОКАЗАТЕЛЬСТВО. Покажем, что при сделанных предположениях выполнено условие (2.4) при  $\tau = \omega$ ,  $B = D + \omega L$ . Действительно,

$$(Bx, x) - \frac{\omega}{2}(Ax, x) = \left(1 - \frac{\omega}{2}\right)(Dx, x) + \frac{\omega}{2}((Lx, x) - (L^T x, x)),$$

но все диагональные элементы положительно определенной матрицы положительны (докажите!), поэтому  $(Dx, x) > 0$  при  $x \neq 0$ , а  $(Lx, x) - (L^T x, x) = (Lx, x) - (x, Lx) = 0 \forall x$ .  $\square$

Естественно, параметр  $\omega$  следует выбирать так, чтобы метод релаксации сходился наиболее быстро. Решение этой задачи далеко выходит за рамки нашего курса. Отметим только, что чаще всего оптимальное значение  $\omega$  лежит вблизи 1,8.

**3. Пример решения задачи оптимизации итерационного параметра.** Рассмотрим итерационный метод (2.3) при  $B = E$ . Этот метод называют методом простой итерации с параметром. Будем предполагать, что  $A$  — симметричная положительно определенная матрица, и выясним, при каких условиях на  $\tau$  метод сходится. Найдем также значение  $\tau$ , обеспечивающее наиболее быструю сходимость.

Напомним нужные в дальнейшем сведения из линейной алгебры. Все характеристические числа симметричной матрицы  $A$  вещественны. Будем нумеровать их в порядке неубывания:  $\lambda_1 \leq \lambda_2, \dots, \leq \lambda_n$ . Существует полная ортонормированная система собственных векторов  $e^1, e^2, \dots, e^n$  матрицы  $A$ :

$$Ae_k = \lambda_k e_k, \quad k = 1, 2, \dots, n, \quad (e_k, e_l) = \begin{cases} 0, & k \neq l; \\ 1, & k = l. \end{cases}$$

Очевидно, что  $\lambda_k = (Ae_k, e_k)$ , поэтому для положительно определенной матрицы  $\lambda_1 > 0$ .

Представим произвольный вектор  $x$  в виде разложения по орто-



нормированному базису собственных векторов матрицы  $A$ :

$$x = \sum_{i=1}^n c_i e_i. \quad (3.1)$$

Тогда  $|x|^2 = \sum_{i=1}^n c_i^2$ ,  $(Ax, x) = \sum_{i=1}^n \lambda_i c_i^2$ , следовательно,

$$\lambda_1 |x|^2 \leq (Ax, x) \leq \lambda_n |x|^2 \quad \forall x. \quad (3.2)$$

Справедлива

**Теорема 3.1.** Пусть  $B = E$ , матрица  $A$  симметрична и положительно определена. Тогда итерационный метод (2.3) сходится при любом начальном приближении, если  $\tau \in (0, 2/\lambda_n)$ . Наилучшая оценка скорости сходимости достигается при  $\tau = \tau_0 = 2/(\lambda_1 + \lambda_n)$ . В этом случае  $|z^k| \leq \rho_0^k |z^0|$ , где  $\rho_0 = (1 - \xi)/(1 + \xi)$ ,  $\xi = \lambda_1/\lambda_n$ .

ДОКАЗАТЕЛЬСТВО. Запишем уравнение для погрешности метода:

$$\frac{z^{k+1} - z^k}{\tau} + Az^k = 0, \quad k = 0, 1, \dots$$

откуда  $z^{k+1} = (E - \tau A)z^k$ . Представим вектор  $z^k$  в виде разложения (3.1). Тогда

$$z^{k+1} = \sum_{i=1}^n (1 - \tau \lambda_i) c_i e^i.$$

Вследствие ортонормированности векторов  $e_1, \dots, e_n$  имеем

$$|z^{k+1}|^2 = \sum_{i=1}^n (1 - \tau \lambda_i)^2 c_i^2, \quad |z^k|^2 = \sum_{i=1}^n c_i^2,$$

откуда

$$|z^{k+1}| \leq \max_{1 \leq i \leq n} |1 - \tau \lambda_i| |z^k|.$$

Понятно, что  $\max_{1 \leq i \leq n} |1 - \tau \lambda_i| \leq \max_{\lambda_1 \leq \lambda \leq \lambda_n} |1 - \tau \lambda|$ , причем поскольку функция  $1 - \tau \lambda$  линейна по  $\lambda$ , то

$$\max_{\lambda_1 \leq \lambda \leq \lambda_n} |1 - \tau \lambda| = \max(|1 - \tau \lambda_1|, |1 - \tau \lambda_n|).$$

Анализируя график функции  $\varphi(\tau) = \max(|1 - \tau \lambda_1|, |1 - \tau \lambda_n|)$ , нетрудно убедиться, что  $0 < \varphi(\tau) < 1$  при  $0 < \tau < 2/\lambda_n$ , а

$$\min_{\tau} \varphi(\tau) = \varphi(\tau_0) = \rho_0. \quad \square \quad (3.3)$$

**4. Итерационные методы вариационного типа.** Отыскание оптимального значения параметра  $\tau$  в рассмотренном выше методе требует предварительного вычисления минимального и максимального собственных чисел матрицы  $A$ . Существуют итерационные методы, позволяющие за счет некоторой дополнительной работы на каждом шаге итераций автоматически настраиваться на оптимальную скорость сходимости. К их числу относятся методы, основанные на замене системы (1.1) эквивалентной задачей минимизации некоторого функционала. В дальнейшем существенную роль играет

**Теорема 4.1.** Пусть матрица  $A$  симметрична и положительно определена. Тогда задача (1.1) эквивалентна задаче отыскания минимума квадратичного функционала  $F(x) = (Ax, x) - 2(b, x)$ .

**ДОКАЗАТЕЛЬСТВО.** Пусть  $x^*$  — решение задачи (1.1), то есть  $Ax^* = b$ . Используя симметрию матрицы  $A$ , получим

$$\begin{aligned} F(x) &= (Ax, x) - 2(Ax^*, x) + (Ax^*, x^*) - (Ax^*, x^*) = \\ &= (A(x - x^*), (x - x^*)) - (Ax^*, x^*), \end{aligned} \quad (4.1)$$

откуда вследствие положительной определенности матрицы  $A$  вытекает, что единственной точкой минимума функционала  $F$  является  $x^*$ .  $\square$

Различные методы минимизации функционала  $F(x)$  приводят к различным итерационным процессам для уравнения (1.1).

Рассмотрим сначала метод покоординатного спуска. Выберем некоторое начальное приближение  $x^0 = (x_1^0, \dots, x_n^0)$  и найдем аргумент  $x_1^1$ , доставляющий минимальное значение функции одной переменной  $F(x_1, x_2^0, \dots, x_n^0)$ . Затем рассмотрим функцию одной переменной  $F(x_1^1, x_2, x_3^0, \dots, x_n^0)$  и найдем точку  $x_2^2$ , в которой она достигает минимума. Выполнив  $n$  таких шагов, построим вектор  $(x_1^1, \dots, x_n^1)$ , примем его за начальное приближение и продолжим описанный процесс.

Используя конкретный вид функционала  $F(x)$ , найдем явные формулы для вычисления векторов  $x^k$ ,  $k = 1, 2, \dots$  в полученном итерационном процессе. Компонента  $x_i^{k+1}$  вектора  $x^{k+1}$  разыскивается как точка минимума функции  $F(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i, x_{i+1}^k, \dots, x_n^k)$ . Выпишем необходимое условие экстремума:

$$F'_{x_i}(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i, x_{i+1}^k, \dots, x_n^k) = 0. \quad (4.2)$$

Вычисляя производную функции  $F(x)$  по переменной  $x_i$ , получим:

$$F'_{x_i}(x) = 2 \sum_{j=1}^n a_{ij} x_j - 2b_i, \quad (4.3)$$

следовательно, решая уравнение (1.1) относительно  $x_i$ , будем иметь

$$x_i^{k+1} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} x_j^{k+1} - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} x_j^k + \frac{b_i}{a_{ii}},$$

и это означает, что метод покоординатного спуска для квадратичного функционала совпадает с методом Зейделя.

Метод релаксации также допускает простую геометрическую интерпретацию. При  $0 < \omega < 1$  из точки  $(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i^k, x_{i+1}^k, \dots, x_n^k)$  двигаются в направлении координатной оси  $x_i$ , не доходя до точки минимума функционала  $F$  на этой прямой, а при  $\omega > 1$  проходят несколько дальше, чем точка минимума функционала. Во многих случаях последний способ приводит к ускорению сходимости.

Опишем еще один метод минимизации функционала. Будем двигаться из точки начального приближения  $x^0$  в направлении наибоыстрейшего убывания функционала  $F$ , т. е. следующее приближение разыскиваем так:  $x^1 = x^0 - \tau \operatorname{grad} F(x^0)$ . Формула (4.3) показывает, что  $\operatorname{grad} F(x^0) = 2(Ax^0 - b)$ . Вектор  $r^0 = Ax^0 - b$  принято называть невязкой. Для сокращения записей удобно обозначить  $2\tau$  вновь через  $\tau$ . Таким образом,  $x^1 = x^0 - \tau r^0$ .

Параметр  $\tau$  выберем так, чтобы значение  $F(x^1)$  было минимальным. Проводя элементарные выкладки, получим  $F(x^1) = F(x^0 - \tau r^0) = F(x^0) - 2\tau(r^0, r^0) + \tau^2(Ar^0, r^0)$ , следовательно, минимум  $F(x^1)$  достигается при  $\tau = \tau_* = (r^0, r^0)/(Ar^0, r^0)$ .

Таким образом, мы пришли к следующему итерационному методу

$$x^{k+1} = x^k - \tau_* r^k, \quad r^k = Ax^k - b, \quad \tau_* = \frac{(r^k, r^k)}{(Ar^k, r^k)}, \quad k = 0, 1, \dots \quad (4.4)$$

Метод (4.4) называют методом наискорейшего спуска. По сравнению с методом простой итерации этот метод требует на каждом шаге итераций проведения дополнительной работы по вычислению параметра  $\tau_*$ . Вследствие этого происходит адаптация к оптимальной скорости сходимости.

**Теорема 4.2.** Пусть матрица  $A$  симметрична и положительно определена. Тогда метод (4.4) сходится при любом начальном приближении. Справедлива следующая оценка скорости сходимости:

$$|z^k| \leq \sqrt{\lambda_n/\lambda_1} \rho_0^k |z^0|. \quad (4.5)$$

**ЗАМЕЧАНИЕ 4.1.** Оценка (4.5) показывает, что погрешность метода (4.4) убывает с той же скоростью, что и погрешность метода простой итерации при оптимальном выборе параметра  $\tau$ .

ДОКАЗАТЕЛЬСТВО теоремы 4.2. Используя (4.4) и вспоминая, что

$$F(x^k - \tau_* r^k) = \min_{\tau} F(x^k - \tau r^k),$$

получим:  $F(x^{k+1}) \leq F(x^k - \tau_0 r^k)$ , где  $\tau_0 = 2/(\lambda_1 + \lambda_n)$ , откуда вследствие представления (4.1) вытекает, что

$$(A(x^{k+1} - x^*), x^{k+1} - x^*) \leq (A(x^k - \tau_0 r^k - x^*), x^k - \tau_0 r^k - x^*).$$

Заметим теперь, что

$$x^k - x^* - \tau_0 r^k = z^k - \tau_0(Ax^k - b) = z^k - \tau_0(Ax^k - Ax^*) = (E - \tau_0 A)z^k,$$

следовательно,

$$(Az^{k+1}, z^{k+1}) \leq (A(E - \tau_0 A)z^k, (E - \tau_0 A)z^k).$$

Представляя вектор  $z^k$  в виде разложения по ортонормированному базису собственных векторов матрицы  $A$ , получим:

$$(A(E - \tau_0 A)z^k, (E - \tau_0 A)z^k) = \sum_{i=1}^n \lambda_i (1 - \tau_0 \lambda_i)^2 c_i^2.$$

Как установлено при доказательстве теоремы (3.1),

$$|1 - \tau_0 \lambda_i| \leq \rho_0 < 1, \quad i = 1, 2, \dots, n,$$

следовательно,  $(Az^{k+1}, z^{k+1}) \leq \rho_0^2 \sum_{i=1}^n \lambda_i c_i^2 = \rho_0^2 (Az^k, z^k)$ , откуда, оче-

видно, вытекает неравенство  $(Az^k, z^k) \leq \rho_0^{2k} (Az^0, z^0)$ . Используя те-

перь оценки (3.2), получим:  $\lambda_1 |z^k|^2 \leq \lambda_n \rho_0^{2k} |z^0|^2$ , а это эквивалент-

но (4.5).  $\square$

### § 3. Методы решения алгебраической проблемы собственных значений

Под алгебраической проблемой собственных значений понимают задачу отыскания собственных чисел и собственных векторов матрицы. Различают полную проблему собственных значений, т. е. нахождение всех собственных чисел и собственных векторов, и частичную проблему собственных значений, т. е. отыскание лишь некоторых собственных чисел и соответствующих им собственных векторов.

Понятно, что методы решения частичной проблемы собственных значений должны быть более простыми. Мы рассмотрим примеры методов обоих классов. При этом ограничимся лишь случаем симметричных матриц.

Собственные числа симметричной матрицы  $A$  будем нумеровать в порядке убывания их модулей:

$$|\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_{n-1}| \leq |\lambda_n|.$$

Через  $e^1, e^2, \dots, e^n$  будем обозначать соответствующие ортонормированные собственные векторы.

**1. Метод прямой итерации.** Этот метод предназначен для отыскания максимального по модулю собственного числа матрицы и соответствующего ему собственного вектора.

Выберем некоторое нормированное начальное приближение  $y^0$  и образуем последовательность нормированных векторов  $y^1, y^2, \dots$ :  $x^1 = Ay^0$ ,  $y^1 = x^1/|x^1|$ . Вообще, по  $y^k$  вычисляем  $x^{k+1} = Ay^k$ , а затем нормируем:  $y^{k+1} = x^{k+1}/|x^{k+1}|$ . Строим также последовательность чисел  $\lambda^{(k)} = (Ay^k, y^k)$ ,  $k = 1, 2, \dots$

**Теорема 1.1.** Пусть  $|\lambda_{n-1}| < |\lambda_n|$ ,  $\theta_k$  — угол, образованный векторами  $y^k$  и  $e^n$ , причем  $\theta_0 \neq \pi/2$ . Тогда  $\operatorname{tg} \theta_k \rightarrow 0$ ,  $\lambda^{(k)} \rightarrow \lambda_n$  при  $k \rightarrow \infty$ . Справедливы следующие оценки скорости сходимости:

$$|\operatorname{tg} \theta_k| \leq \left| \frac{\lambda_{n-1}}{\lambda_n} \right|^k |\operatorname{tg} \theta_0|, \quad |\lambda^{(k)} - \lambda_n| \leq (|\lambda_{n-1}| + |\lambda_n|) \left| \frac{\lambda_{n-1}}{\lambda_n} \right|^{2k}. \quad (1.1)$$

**ДОКАЗАТЕЛЬСТВО.** По теореме об ортогональном разложении евклидова пространства вектор  $y^k$  однозначно представим в виде

$$y^k = \cos \theta_k e^n + \sin \theta_k u^k, \quad (1.2)$$

где  $u^k$  — единичный вектор, ортогональный  $e^n$ . В соответствии с изучаемым алгоритмом

$$y^{k+1} = \alpha(\cos \theta_k A e^n + \sin \theta_k A u^k) = \alpha(\cos \theta_k \lambda_n e^n + \sin \theta_k A u^k).$$

Здесь  $\alpha$  — число, выбираемое так, чтобы вектор  $y^{k+1}$  имел единичную длину. Перепишем последнее равенство в виде

$$\begin{aligned} y^{k+1} &= \alpha(\cos \theta_k \lambda_n e^n + |A u^k| \sin \theta_k A u^k / |A u^k|) = \\ &= \cos \theta_{k+1} e^n + \sin \theta_{k+1} A u^k / |A u^k|. \end{aligned} \quad (1.3)$$

Мы учли здесь, что вектор  $Au^k$  ортогонален  $e^n$ , поскольку

$$(Au^k, e^n) = (u^k, Ae^n) = \lambda_n(u^k, e^n) = 0.$$

Из (1.3) вытекает, что

$$\operatorname{tg} \theta_{k+1} = \frac{|Au^k|}{\lambda_n} \operatorname{tg} \theta_k.$$

Записывая разложение  $u^k$  по базису собственных векторов матрицы  $A$  и учитывая ортогональность  $u^k$  и  $e^n$ , получим, что  $Au^k = \sum_{i=1}^{n-1} c_i \lambda_i e^i$ , следовательно,

$$|Au^k|^2 = \sum_{i=1}^{n-1} c_i^2 \lambda_i^2 \leq \lambda_{n-1}^2 \sum_{i=1}^{n-1} c_i^2 = \lambda_{n-1}^2 |u^k|^2 = \lambda_{n-1}^2. \quad (1.4)$$

Таким образом,

$$|\operatorname{tg} \theta_{k+1}| \leq \frac{|\lambda_{n-1}|}{|\lambda_n|} |\operatorname{tg} \theta_k|,$$

и первая оценка (1) доказана.

Вновь используя представление (1.2) и то, что  $(Au^k, e^n) = 0$ , получим:

$$\begin{aligned} \lambda_n - \lambda^{(k)} &= \lambda_n - (A(\cos \theta_k e^n + \sin \theta_k u^k), \cos \theta_k e^n + \sin \theta_k u^k) = \\ &= \lambda_n - (\cos \theta_k \lambda_n e^n + \sin \theta_k Au^k, \cos \theta_k e^n + \sin \theta_k u^k) = \\ &= \lambda_n - (\lambda_n \cos^2 \theta_k + (Au^k, u^k) \sin^2 \theta_k) = (\lambda_n - (Au^k, u^k)) \sin^2 \theta_k, \end{aligned}$$

Откуда вследствие (1.4) вытекает, что

$$|\lambda_n - \lambda^{(k)}| \leq |\lambda_n + \lambda_{n-1}| \sin^2 \theta_k \leq |\lambda_n + \lambda_{n-1}| \operatorname{tg}^2 \theta_k,$$

Вместе с уже полученной первой оценкой (1.1) это завершает доказательство теоремы.  $\square$

Условие  $\theta_0 \neq \pi/2$  на практике не слишком обременительно. Если оно нарушается, то при проведении итераций за счет ошибок округления приближения обязательно выйдут из гиперплоскости, ортогональной  $e^n$ .

**2. Метод обратной итерации.** Метод предназначен для отыскания минимального по модулю собственного числа и соответствующего ему собственного вектора. Опишем соответствующий алгоритм.

Выбираем нормированное начальное приближение  $y^0$  и строим последовательность векторов  $y^1, y^2, \dots$  по формулам:  $x^{k+1} = A^{-1}y^k$ ,  $y^{k+1} = x^{k+1}/|x^{k+1}|$ , а также числа  $\lambda^{(k)} = (Ay^k, y^k)$ ,  $k = 0, 1, 2, \dots$

При реализации метода выгоднее не строить и хранить матрицу  $A^{-1}$ , а решать на каждой итерации систему линейных уравнений  $Ax^{k+1} = y^k$ . Предварительно целесообразно представить матрицу  $A$  в виде  $LU$  или  $QU$  разложения (см. §1).

Относительно сходимости метода справедлива теорема, полностью аналогичная теореме 1.1, но на этот раз скорость сходимости характеризуется отношением  $|\lambda_1|/|\lambda_2| < 1$ .

**2.1. Метод обратной итерации со сдвигом.** Рассмотрим обобщение предыдущего метода, а именно переход от вектора  $y^k$  к  $y^{k+1}$  будем выполнять по формулам:  $(A - \sigma E)x^{k+1} = y^k$ ,  $y^{k+1} = x^{k+1}/|x^{k+1}|$ . Здесь  $\sigma$  — параметр, называемый сдвигом. Последовательность  $\lambda^{(k)}$ ,  $k = 1, 2, \dots$ , по-прежнему, определяется формулой  $\lambda^{(k)} = (Ay^k, y^k)$ . Сходимость этого метода исследуется по той же схеме, что и в теореме 1.1. При этом оказывается, что  $\lambda^{(k)} \rightarrow \lambda_j$ , где номер  $j$  характеризуется условием  $|\lambda_j - \sigma| < |\lambda_i - \sigma| \quad \forall i \neq j$ , а последовательность  $y^k$  сходится к соответствующему собственному вектору  $e^j$ . Таким образом, метод позволяет находить собственное число матрицы  $A$ , ближайшее к заданному числу  $\sigma$ .

**3. Метод вращений (Якоби).** Этот метод довольно часто используется при решении полной проблемы собственных значений для симметричных матриц не слишком высокого порядка.

Напомним определение матрицы вращения. Ортогональная матрица

$$\begin{pmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{pmatrix}$$

порождает преобразование поворота на угол  $\varphi$  в двумерной плоскости. Матрица  $Q = (q_{ij})_{i,j=1}^n$ , отличающаяся от единичной лишь четырьмя элементами:  $q_{k,k} = \cos \varphi$ ,  $q_{ll} = \cos \varphi$ ,  $q_{kl} = \sin \varphi$ ,  $q_{lk} = -\sin \varphi$ , где  $1 \leq k < l \leq n$  — заданные целые числа, называется матрицей вращения. Нетрудно проверить, что  $Q$  — ортогональная матрица. Она порождает преобразование поворота на угол  $\varphi$  в двумерной плоскости, натянутой на векторы канонического базиса с номерами  $k, l$ .

Опишем идею метода Якоби. Пусть  $A$  — симметричная матрица. Как и в §1, ее собственные числа будем нумеровать в порядке неубывания. Образует матрицу  $T$ , столбцами которой являются ортонормированные собственные векторы  $e^1, e^2, \dots, e^n$  матрицы  $A$ . Нетрудно

убедиться, что  $T^T AT = \Lambda$ , где  $\Lambda$  — диагональная матрица с элементами  $\lambda_1, \lambda_2, \dots, \lambda_n$  на диагонали. В методе Якоби матрица  $T$  строится как предел последовательности ортогональных матриц  $T_s$  так, что  $\lim_{s \rightarrow \infty} T_s^T AT_s = \Lambda$ , причем при каждом  $s$  матрица  $T_s$  конструируется как произведение матриц вращения.

Образую по матрице  $A$  матрицу  $\hat{A} = Q^T A Q$  и попытаемся выбрать параметры матрицы вращения, то есть  $k, l, \varphi$  так, чтобы матрица  $\hat{A}$  была максимально близка к диагональной. Опуская элементарные выкладки, приведем выражение для суммы квадратов внедиагональных элементов матрицы  $\hat{A}$ :

$$\sum_{i \neq j} \hat{a}_{ij}^2 = \sum_{i \neq j} a_{ij}^2 - 2a_{kl}^2 + 2[a_{kl} \cos 2\varphi + \frac{1}{2}(a_{ll} - a_{kk}) \sin 2\varphi]^2.$$

Определим теперь числа  $k, l$  из условия:

$$|a_{kl}| = \max_{i \neq j} |a_{ij}| \quad (3.1)$$

а затем угол  $\varphi$  так, чтобы

$$a_{kl} \cos 2\varphi + \frac{1}{2}(a_{ll} - a_{kk}) \sin 2\varphi = 0,$$

или

$$\operatorname{tg} 2\varphi = \frac{2a_{kl}}{a_{kk} - a_{ll}}. \quad (3.2)$$

При указанном выборе параметров матрицы вращения сумма квадратов внедиагональных элементов матрицы  $\hat{A}$  принимает наименьшее значение.

Теперь можно описать алгоритм метода Якоби. Пусть  $A_0 = A$ . Образую последовательность матриц  $A_1, A_2, \dots$ , при помощи рекуррентной формулы:

$$A_{p+1} = Q_p^T A_p Q_p, \quad p = 0, 1, \dots, \quad (3.3)$$

где параметры матрицы вращения  $Q_p$  определяются так, чтобы сделать сумму квадратов внедиагональных элементов матрицы  $A_{p+1}$  минимально возможной, т. е. по формулам вида (3.1), (3.2).

Итерации проводят до тех пор, пока все внедиагональные элементы матрицы  $A_p$  не станут достаточно малыми. Тогда в качестве приближений к собственным числам матрицы  $A$  принимают диагональные элементы матрицы  $A_p$ , а столбцы матрицы  $Q_0 Q_1 \dots Q_p$  считают приближениями к собственным векторам матрицы  $A$ .

При обосновании сходимости метода используется



**Лемма 3.1.** Пусть параметры матрицы вращения  $Q$  определяются согласно формулам (3.1), (3.2). Тогда

$$\sum_{i \neq j} \hat{a}_{ij}^2 \leq q \sum_{i \neq j} a_{ij}^2, \quad (3.4)$$

где

$$0 < q = 1 - \frac{2}{n(n-1)} < 1$$

при  $n > 2$ .

ДОКАЗАТЕЛЬСТВО. Вследствие (3.2) имеем:

$$\sum_{i \neq j} \hat{a}_{ij}^2 = \sum_{i \neq j} a_{ij}^2 - 2a_{kl}^2, \quad (3.5)$$

а на основании (3.1)

$$\sum_{i \neq j} a_{ij}^2 \leq a_{kl}^2 n(n-1). \quad (3.6)$$

Здесь учтено, что матрица порядка  $n$  имеет  $n^2 - n$  внедиагональных элементов. Из (3.5), (3.6) очевидным образом следует (3.4).  $\square$

Будем опираться также на следующий фундаментальный результат теории симметричных матриц, который приведем без доказательства.

**Теорема 3.1.** Пусть  $A, B$  — симметричные матрицы порядка  $n$ ,  $\lambda_k(A)$ ,  $\lambda_k(B)$ ,  $k = 1, \dots, n$ , — их собственные числа, занумерованные в порядке убывания. Тогда

$$(\lambda_k(A) - \lambda_k(B))^2 \leq \sum_{i,j=1}^n (a_{ij} - b_{ij})^2, \quad k = 1, 2, \dots, n. \quad (3.7)$$

Докажем теперь сходимость метода Якоби. Из рекуррентной формулы (3.3) и леммы ?? вытекает, что

$$\sum_{i \neq j} (a_{ij}^{(p)})^2 \leq q^p \sum_{i \neq j} a_{ij}^2 \rightarrow 0 \text{ при } p \rightarrow \infty.$$

Это значит, что по любому заданному  $\varepsilon > 0$  можно указать целое положительное число  $p$  такое, что

$$\sum_{i \neq j} (a_{ij}^{(p)})^2 \leq \varepsilon^2. \quad (3.8)$$

Обозначим через  $\Lambda_p$  диагональную матрицу, на диагонали которой расположены диагональные элементы матрицы  $A_p$ . В соответствии с теоремой 3.1 и оценкой (3.8) можем написать:

$$|\lambda_k(A_p) - \lambda_k^{(p)}| \leq \varepsilon, \quad k = 1, 2, \dots, n,$$

где  $\lambda_k^{(p)}$ ,  $k = 1, \dots, n$ , — диагональные элементы матрицы  $\Lambda_p$ , упорядоченные по неубыванию. Вследствие (3.3) имеем:  $A_p = T_p^T A T_p$ , где  $T_p = Q_0 Q_1 \dots Q_p$ , т. е. матрицы  $A_p$  и  $A$  подобны, а, значит, их собственные числа совпадают, поэтому

$$|\lambda_k(A) - \lambda_k^{(p)}| \leq \varepsilon, \quad k = 1, 2, \dots, n. \quad (3.9)$$

Сходимость метода Якоби доказана.

Подчеркнем, что оценка (3.9) может рассматриваться как апостериорная оценка погрешности метода: добившись в ходе итераций выполнения неравенства (3.8), мы получаем собственные числа с погрешностью, не превосходящей  $\varepsilon$ .

## § 4. Методы решения нелинейных уравнений

В этом параграфе излагаются наиболее распространенные методы решения нелинейных уравнений вида

$$f(x) = 0, \quad (1)$$

где  $f$  — заданная непрерывная функция вещественного переменного. Всюду в дальнейшем предполагается, что известен отрезок  $[x_0, x_1]$ , содержащий корень  $\alpha$  уравнения (1). Предполагается, что  $\alpha$  — единственный корень уравнения (1) на отрезке  $[x_0, x_1]$ . Методы отделения корней существенно зависят от конкретного вида функции  $f$  и в настоящем пособии не рассматриваются.

**1. Метод деления отрезка пополам.** Будем считать, что  $f(x_0)f(x_1) < 0$ , т. е. функция  $f$  меняет знак в точке  $\alpha \in [x_0, x_1]$ . Положим  $x_2 = (x_0 + x_1)/2$  и вычислим  $f(x_0)f(x_2)$ . Если  $f(x_0)f(x_2) < 0$ , то корень расположен на отрезке  $[x_0, x_2]$ . В противном случае — на отрезке  $[x_2, x_1]$ . Выбирая теперь тот из двух отрезков, на котором лежит корень уравнения, применяем к нему описанную процедуру деления пополам. Процесс продолжают до тех пор, пока длина отрезка не станет меньше заданного  $\varepsilon > 0$ . При этом корень, очевидно, будет найден с точностью  $\varepsilon > 0$ .

Метод привлекает своей простотой и малыми затратами для реализации каждого тога итерационного процесса (требуется лишь вычисления одного значения функции  $f$ ). Следует подчеркнуть, что метод неприменим для отыскания корней четной кратности, когда функция  $f$  обращается в нуль в точке  $\alpha \in [x_0, x_1]$ , но не меняет знака на отрезке  $[x_0, x_1]$ .

**2. Метод простой итерации.** Метод основан на приведении уравнения (1) к эквивалентному уравнению вида

$$x = \varphi(x) \quad (2.1)$$

и построении последовательности приближений  $x_0, x_1, \dots, x_n, \dots$  по формуле

$$x_{k+1} = \varphi(x_k), \quad k = 0, 1, 2, \dots \quad (2.2)$$

Начальное приближение  $x_0$  считается заданным.

Эквивалентное преобразование уравнения (1) к виду (2.1) может быть выполнено различными способами. Например, можно положить  $\varphi(x) = x + \psi(x)f(x)$ , где  $\psi$  — произвольная непрерывная функция, не обращающаяся в нуль в окрестности точки  $\alpha$ . Функцию  $\psi$  следует конструировать так, чтобы обеспечить сходимость последовательности  $x_k$ . Тогда, очевидно,  $\lim_{k \rightarrow \infty} x_k = \alpha$ .

**Теорема 2.1.** Пусть в некоторой окрестности

$$S(\alpha, r) = \{x, |x - \alpha| < r\}$$

корня уравнения (2.1) функция  $\varphi$  удовлетворяет условию Липшица с постоянной меньшей единицы, т. е.

$$|\varphi(x') - \varphi(x'')| \leq q|x' - x''|, \quad x', x'' \in S(\alpha, r), \quad 0 < q < 1. \quad (2.3)$$

Тогда при любом  $x_0 \in S(\alpha, r)$

$$|x_k - \alpha| \leq q^k |x_0 - \alpha| \rightarrow 0 \quad \text{при } k \rightarrow \infty. \quad (2.4)$$

**ДОКАЗАТЕЛЬСТВО.** Если  $x_0 \in S(\alpha, r)$ , то

$$|x_1 - \alpha| = |\varphi(x_0) - \alpha| = |\varphi(x_0) - \varphi(\alpha)| \leq q|x_0 - \alpha| < r,$$

следовательно,  $x_1 \in S(\alpha, r)$ . Аналогично,  $x_1, x_2, \dots \in S(\alpha, r)$ . Таким образом,

$$|x_k - \alpha| = |\varphi(x_{k-1}) - \varphi(\alpha)| \leq q|x_{k-1} - \alpha| \leq q^k |x_0 - \alpha| \rightarrow 0, \quad k \rightarrow \infty. \quad \square$$

**Следствие 2.1.** Пусть функция  $\varphi$  непрерывно дифференцируема на  $S(\alpha, r)$  и

$$|\varphi'(x)| \leq q < 1 \quad (2.5)$$

при  $x \in S(\alpha, r)$ . Тогда для последовательности  $x_k$ , построенной по методу (2.2), выполняется оценка (2.4).

Графическая интерпретация итерационного процесса (2.2) дана на рис. 1, 2. Отметим, что случай  $-1 < \varphi'(x) \leq 0$  (рис. 2) наиболее

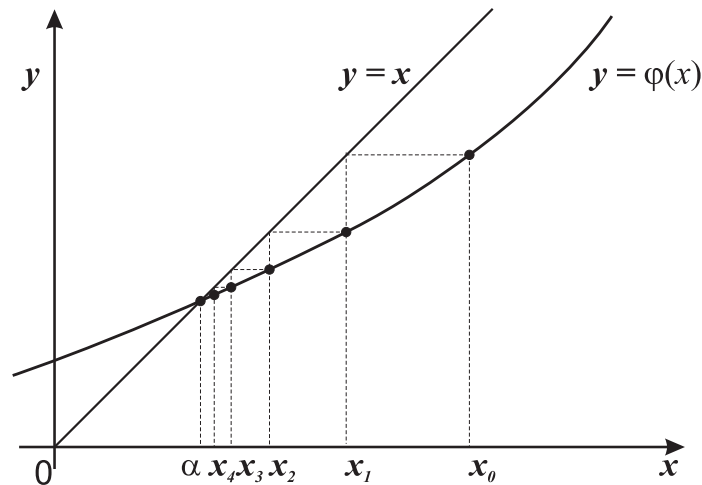


Рис. 1. К итерационному процессу (2.2),  $0 \leq \varphi'(x) < 1$

выгоден, так как величина  $|x_k - x_{k-1}|$  на каждом шаге итерационного процесса оценивает сверху погрешность приближения  $x_k$  к  $\alpha$ .

Приведем простейший пример конструирования функции  $\varphi$ . Пусть требуется вычислить  $\sqrt{a}$ ,  $a > 0$ . Соответствующее уравнение имеет вид

$$f(x) \equiv x^2 - a = 0. \quad (2.6)$$

Преобразуем уравнение (2.6):

$$x = \varphi(x), \quad \text{где } \varphi(x) = \frac{a}{x}.$$

В этом случае  $\varphi'(x) = -a/x^2$ ,  $\varphi'(\sqrt{a}) = -1$ , условие (2.5) в окрестности корня не выполнено, и итерационный процесс может разойтись, даже если  $x_0$  лежит сколь угодно близко к  $\sqrt{a}$  (проиллюстрируйте эту ситуацию графически!).

Перепишем теперь уравнение (2.6) в виде

$$x = \frac{1}{2} \left( x + \frac{a}{x} \right)$$

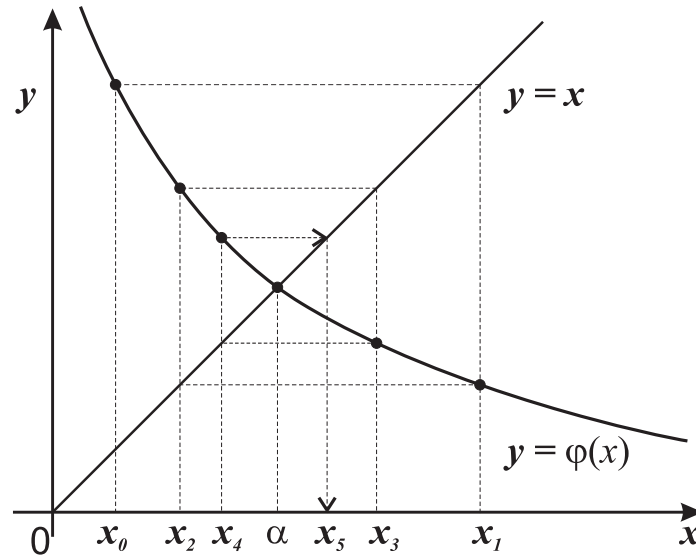


Рис. 2. К итерационному процессу (2.2),  $-1 < \varphi'(x) \leq 0$

В этом случае

$$\varphi(x) = \frac{1}{2} \left( x + \frac{a}{x} \right), \quad \varphi'(x) = \frac{1}{2} \left( 1 - \frac{a}{x^2} \right), \quad \varphi'(\sqrt{a}) = 0$$

и, следовательно,  $|\varphi'(x)| < 1$  в некоторой окрестности  $\sqrt{a}$ . Таким образом, итерационный метод, определяемый соотношением

$$x_{k+1} = \frac{1}{2} \left( x_k + \frac{a}{x_k} \right), \quad k = 0, 1, \dots,$$

сходится при любом начальном приближении  $x_0$ , достаточно близком к  $\sqrt{a}$ . Более детальный анализ показывает, что сходимость имеет место при любом  $x_0 > 0$ .

**2.1.** Заметим, что если

$$\varphi'(\alpha) = \varphi''(\alpha) = \dots = \varphi^{(m-1)}(\alpha) = 0, \quad \varphi^{(m)}(\alpha) \neq 0, \quad \text{при } m \geq 2,$$

то применяя формулу Тейлора с остаточным членом в форме Лагранжа, получим, что

$$|x_{k+1} - \alpha| \leq \frac{1}{m!} |\varphi^{(m)}(\xi)| |x_k - \alpha|^m, \quad (2.7)$$

и метод (2.2) сходится вблизи корня очень быстро. Число  $m$  принято называть порядком итерационного метода (2.2).

**3. Метод Ньютона.** В этом случае полагают

$$\varphi(x) = x - \frac{f(x)}{f'(x)},$$

то есть

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, \dots \quad (3.1)$$

Нетрудно видеть, что в этом случае

$$\varphi'(\alpha) = 0, \quad \varphi''(\alpha) \neq 0, \quad (3.2)$$

если  $f(\alpha) = 0$ , а  $f'(\alpha) \neq 0$ , т. е.  $\alpha$  — простой корень уравнения (1). Отсюда вытекает, что метод Ньютона (3.1) сходится при любом начальном приближении  $x_0$ , достаточно близком к  $\alpha$ .

Равенства (3.2) показывают, что метод Ньютона — метод второго порядка и, следовательно, имеет квадратичную сходимость (в оценке (2.7) получаем  $m = 2$ ).

Метод Ньютона имеет простую геометрическую интерпретацию:  $x_{k+1}$  — точка пересечения касательной к графику функции  $f$  в точке  $(x_k, f(x_k))$  с осью  $x$  (рис. 3).

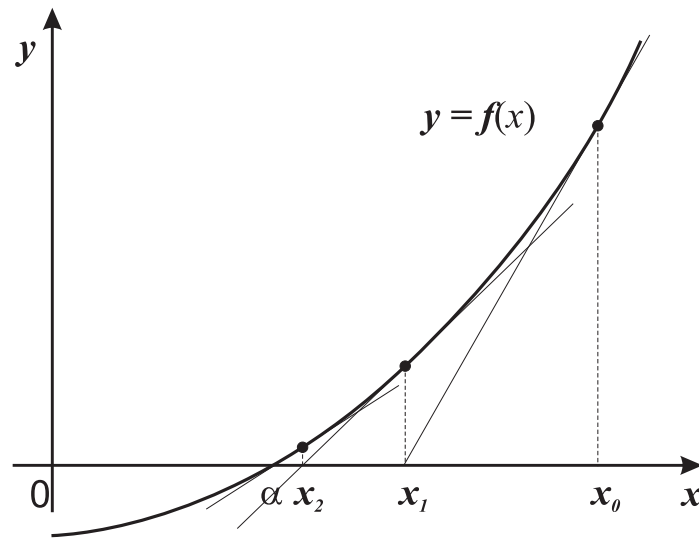


Рис. 3. Метод Ньютона

**4. Метод хорд.** Недостатком метода Ньютона является необходимость вычисления производной функции  $f$ . В некоторых случаях это оказывается слишком трудоемкой задачей.

Если известен интервал  $(x_0, x_1)$ , на концах которого функция  $f$  имеет противоположные знаки, то приближение  $x_2$  к корню можно

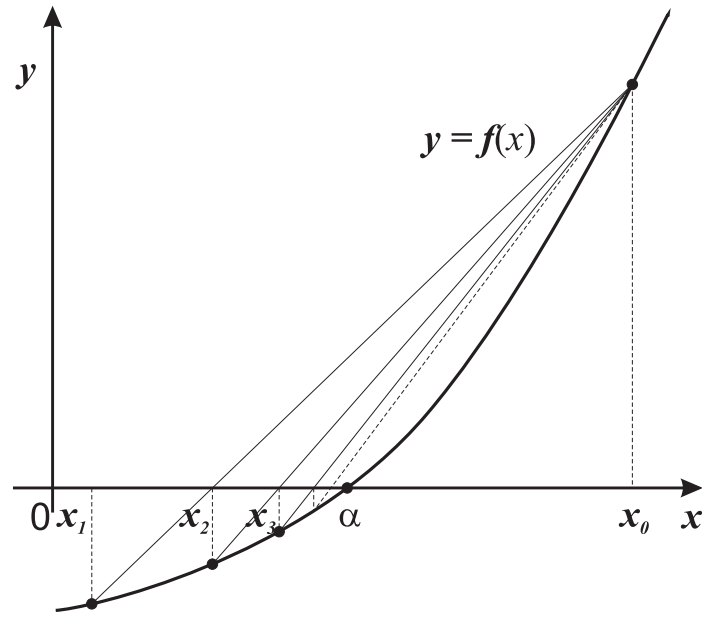


Рис. 4. Метод хорд

определить как точку пересечения прямой, проходящей через точки  $(x_0, f(x_0))$  и  $(x_1, f(x_1))$ , с осью  $x$  (рис. 4). Ясно, что

$$x_2 = \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)}.$$

Все последующие приближения будем вычислять по формуле

$$x_{k+1} = \frac{x_0 f(x_k) - x_k f(x_0)}{f(x_k) - f(x_0)}, \quad k = 2, 3, \dots \quad (4.1)$$

Метод (4.1) можно переписать в виде (2.2), где

$$\varphi(x) = \frac{x_0 f(x) - x f(x_0)}{f(x) - f(x_0)}.$$

Используя формулу Тейлора, нетрудно подсчитать, что

$$\varphi'(\alpha) = \frac{(x_0 - \alpha)^2 f''(\xi)}{2 f'(\xi_1)}, \quad \xi, \xi_1 \in (x_0, x_1),$$

откуда вытекает, что  $|\varphi'(\alpha)| < 1$ , если  $x_0$  выбрано достаточно близко к корню. Таким образом, существует окрестность точки  $\alpha$ , в которой  $|\varphi'(x)| \leq q < 1$ , и метод (4.1) сходится при любом  $x_1$  из этой окрестности.

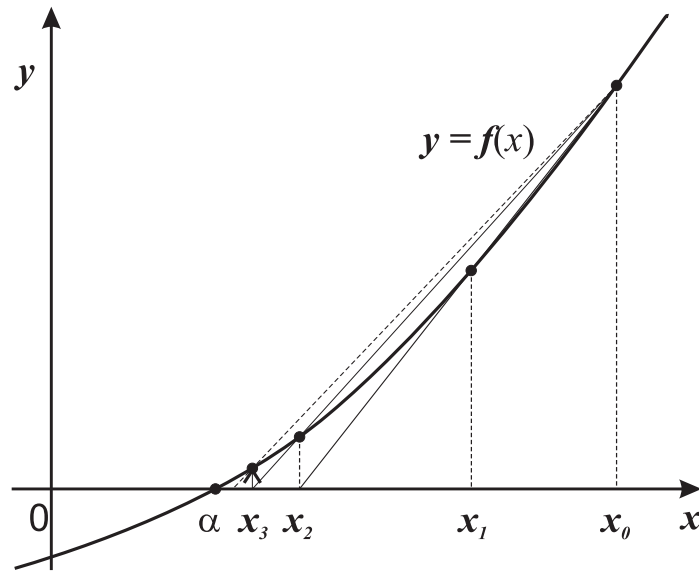


Рис. 5. Метод секущих

**5. Метод секущих.** Этот метод получается при аппроксимации производной  $f'(x_k)$  в методе Ньютона (3.1) разностным отношением

$$f'(x_k) = \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}.$$

В результате приходим к следующему итерационному процессу

$$x_{k+1} = x_k - \frac{(x_k - x_{k-1})f(x_k)}{f(x_k) - f(x_{k-1})} \quad (5.1)$$

(рис. 5). Геометрически  $x_{k+1}$  — точка пересечения секущей, проходящей через точки  $(x_k, f(x_k))$  и  $(x_{k-1}, f(x_{k-1}))$  с осью  $x$ .

В отличие от ранее рассмотренных методов для вычисления каждого последующего приближения  $x_{k+1}$  требуется использовать два предыдущих приближения. Такие методы принято называть двушаговыми.

Вновь используя разложение по формуле Тейлора, можно показать, что справедливо приближенное равенство

$$x_{k+1} - \alpha \approx \frac{1}{2} \frac{f''(\alpha)}{f'(\alpha)} (x_{k-1} - \alpha)(x_k - \alpha),$$

показывающее, что если имеет место сходимость метода (5.1), то его погрешность убывает существенно быстрее, чем в методе хорд, и несколько медленнее, чем в методе Ньютона.





**2. Метод Ньютона.** При наличии достаточно хорошего приближения к решению системы (2) очень часто весьма эффективным оказывается его уточнение по методу Ньютона. Этот метод является непосредственным обобщением метода Ньютона для одного нелинейного уравнения и формально может быть записан следующим образом

$$x^{k+1} = x^k - (F'(x^k))^{-1}F(x^k), \quad k = 0, 1, 2, \dots \quad (2.1)$$

Здесь

$$F'(x^k) = \left\{ \frac{\partial f_i(x^k)}{\partial x_j} \right\}_{i,j=1}^n$$

есть матрица Якоби системы функций  $f_1, f_2, \dots, f_n$ , вычисленная в точке  $x^k$ .

При практической реализации метода его обычно переписывают в виде

$$F'(x^k)(x^{k+1} - x^k) = -F(x^k),$$

или

$$\begin{aligned} F'(x^k)\Delta^k &= -F(x^k), \\ x^{k+1} &= x^k + \Delta^k. \end{aligned}$$

Таким образом, для построения  $x^{k+1}$  по известному  $x^k$  нужно решить систему линейных алгебраических уравнений с матрицей  $F'(x^k)$ .

Во многих случаях формирование матрицы Якоби  $F'(x^k)$  оказывается существенно более трудоемкой задачей, чем вычисление вектора  $F(x^k)$ . Поэтому довольно часто применяют следующую модификацию метода Ньютона

$$\begin{aligned} F'(x^0)\Delta^k &= -F(x^k), \\ x^{k+1} &= x^k + \Delta^k. \end{aligned}$$

Здесь матрица Якоби вычисляется лишь на начальном приближении. Экономия достигается как за счет формирования матрицы системы, так и за счет того, что решение нескольких систем с одинаковыми матрицами существенно проще, чем решение того же количества систем с разными матрицами.

При этом надо, однако, иметь в виду, что если метод (2.1) обладает при достаточно хорошем начальном приближении квадратичной сходимостью, т. е.

$$|x^{k+1} - \alpha| \leq \text{const}|x^k - \alpha|^2,$$

то модифицированный метод Ньютона сходится лишь со скоростью геометрической прогрессии.

---

---

ГЛАВА 2  
**Численные методы анализа**

**§ 1. Интерполирование функций**

**1. Существование и единственность интерполяционного полинома.** Под интерполированием понимают следующую задачу. Даны  $n + 1$  различные точки  $x_0, x_1, \dots, x_n$  на вещественной прямой и значения некоторой вещественной функции  $f(x_0), f(x_1), \dots, f(x_n)$  в этих точках. Даны также вещественные функции  $\varphi_0(x), \dots, \varphi_n(x)$ . Требуется построить функцию

$$\varphi(x) = \sum_{k=0}^n c_k \varphi_k(x), \quad (1.1)$$

где  $c_k$  — искомые вещественные числа, удовлетворяющую условиям:

$$\varphi(x_i) = f(x_i), \quad i = 0, 1, \dots, n. \quad (1.2)$$

Функцию  $\varphi$  принято называть обобщенным интерполяционным полиномом, функции  $\varphi_k$  — базисными функциями, точки  $x_0, x_1, \dots, x_n$  называют узлами интерполирования. Если  $\varphi_k$  — полиномы, то говорят о задаче алгебраического интерполирования, а функцию  $\varphi$  называют интерполяционным полиномом. Задачей алгебраического интерполирования мы в дальнейшем и будем заниматься.

Если  $\varphi_0(x) \equiv 1, \varphi_1(x) = x, \varphi_2(x) = x^2, \dots, \varphi_n(x) = x^n$ , то  $\varphi(x)$  есть полином степени  $n$ .

**Теорема 1.1.** *Интерполяционный полином степени  $n$  однозначно определяется по любым заданным значениям  $f(x_0), f(x_1), \dots, f(x_n)$  в  $n + 1$  различных точках  $x_0, x_1, \dots, x_n$ .*

**ДОКАЗАТЕЛЬСТВО.** Условия (1.2) представляют собой систему из  $n + 1$  линейных алгебраических уравнений относительно  $n + 1$  неизвестных  $c_0, c_1, \dots, c_n$ , поэтому достаточно установить, что соответствующая однородная система имеет только тривиальное решение. Предполагая противное, получим, что полином степени  $n$  имеет  $n + 1$  различных корней  $x_0, x_1, \dots, x_n$ , а это возможно лишь в том случае, когда все его коэффициенты, т. е.  $c_0, c_1, \dots, c_n$ , равны нулю.  $\square$

Во многих случаях полезно иметь явный вид интерполяционного многочлена. Однако разрешить систему (1.2) в явном виде при

произвольном значении  $n$  не удастся. Причина — в неудачном выборе формы представления интерполяционного полинома, иными словами, в неудачном выборе базисных функций  $\varphi_0(x), \dots, \varphi_n(x)$ .

**2. Интерполяционный полином Лагранжа.** Понятно, что наиболее простой вид система (1.2) принимает в том случае, когда базисные функции — полиномы степени  $n$ , удовлетворяющие условиям:

$$\varphi_i(x_j) = \begin{cases} 0, & i \neq j; \\ 1, & i = j. \end{cases} \quad (2.1)$$

При этом, очевидно,  $c_i = f(x_i)$ ,  $i = 0, 1, \dots, n$ . Базисные функции в этом случае легко выписываются в явном виде. Действительно, вследствие условия (2.1) точки  $x_0, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$  — корни полинома  $\varphi_i(x)$ , поэтому  $\varphi_i(x) = A(x-x_0) \cdots (x-x_{i-1})(x-x_{i+1}) \cdots (x-x_n)$ , где  $A = \text{const}$ . Используя условие  $\varphi_i(x_i) = 1$ , находим

$$A = \frac{1}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)},$$

следовательно,

$$\varphi_i(x) = \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)}, \quad (2.2)$$

$i = 0, \dots, n$ . Построенный таким образом интерполяционный полином называют интерполяционным полиномом Лагранжа и обозначают через  $L_n(x)$ :

$$L_n(x) = \sum_{i=0}^n f(x_i) \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)}. \quad (2.3)$$

Полиномы, определяемые формулами (2.2), называют базисными функциями Лагранжа.

**3. Интерполяционный полином Ньютона.** Определим теперь базисные функции соотношениями:

$$\begin{aligned} \varphi_0(x) &\equiv 1, \quad \varphi_1(x) = (x - x_0), \quad \varphi_2(x) = (x - x_0)(x - x_1), \dots, \\ \varphi_n(x) &= (x - x_0)(x - x_1) \cdots (x - x_{n-1}). \end{aligned} \quad (3.1)$$

Тогда система (1.2) принимает треугольный вид:

$$c_0 = f(x_0),$$



Величину  $f(x_0, x_1, \dots, x_k)$  называют разделенной разностью порядка  $k$ .

Таким образом, интерполяционный полином приобретает вид

$$P_n(x) = f(x_0) + f(x_0, x_1)(x - x_0) + f(x_0, x_1, x_2)(x - x_0)(x - x_1) + \dots + f(x_0, x_1, x_2, \dots, x_n)(x - x_0)(x - x_1) \cdots (x - x_{n-1}). \quad (3.5)$$

Такую форму записи интерполяционного полинома называют интерполяционным полиномом Ньютона.

Процесс вычисления разделенных разностей полезно иллюстрировать таблицей.

$x_0$	$f(x_0)$				
	$f(x_1)$	$\frac{f(x_0, x_1)}{f(x_1, x_2)}$			
	$f(x_2)$	$\frac{f(x_1, x_2)}{f(x_2, x_3)}$	$\dots$	$\frac{f(x_0, x_1, x_2, \dots, x_n)}{f(x_1, x_2, x_3)}$	
$\vdots$	$\vdots$				
$x_n$	$f(x_n)$				

Первые два столбца этой таблицы — исходные данные. Остальные — вычисляются последовательно один за другим. Подчеркнутые разделенные разности (они стоят на верхней диагонали таблицы) используются при образовании полинома  $P_n(x)$ .

#### 4. Оценка погрешности интерполирования.

**Теорема 4.1.** Пусть  $[a, b]$  — отрезок вещественной оси, содержащий все узлы интерполирования  $x_0, x_1, \dots, x_n$ , функция  $f$  непрерывно дифференцируема  $n + 1$  раз на отрезке  $[a, b]$ , точка  $x \in [a, b]$  и не совпадает ни с одним из узлов интерполирования. Тогда существует точка  $\xi \in [a, b]$  такая, что<sup>1)</sup>

$$f(x) = L_n(x) + \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{n+1}(x). \quad (4.1)$$

Здесь  $L_n$  — полином степени  $n$ , интерполирующий функцию  $f$  по узлам  $x_0, x_1, \dots, x_n$ ,  $\omega_{n+1}(x) = (x - x_0)(x - x_1) \cdots (x - x_n)$ .

**ДОКАЗАТЕЛЬСТВО.** Введем в рассмотрение функцию

$$g(z) = f(z) - L_n(z) - K\omega_{n+1}(z), \quad z \in [a, b].$$

<sup>1)</sup>Выражение в правой части (4.1) принято называть остаточным членом интерполяционного полинома в форме Лагранжа.

Здесь  $K$  — постоянная. Определим ее так, чтобы

$$f(x) - L_n(x) - K\omega_{n+1}(x) = 0. \quad (4.2)$$

Это можно сделать, так как  $\omega_{n+1}(x) \neq 0$ . При указанном выборе постоянной  $K$  функция  $g$  обращается в нуль в  $n + 2$  точках, а именно в точке  $x$ , а также в точках  $x_0, x_1, \dots, x_n$ , так как  $L_n(x_i) - f(x_i) = 0$ ,  $\omega_{n+1}(x_i) = 0$  при  $i = 0, 1, \dots, n$ . По теореме Ролля первая производная функции  $g$  обращается в нуль по крайней мере один раз между каждой парой соседних точек из указанного множества точек. Иными словами, существует  $n + 1$  различных точек  $\xi_1, \xi_2, \dots, \xi_{n+1} \in [a, b]$  таких, что  $g'(\xi_i) = 0$ ,  $i = 1, \dots, n + 1$ . Отсюда вытекает, что вторая производная функции  $g$  обращается в нуль по крайней мере в  $n$  различных точках. Продолжая аналогичные рассуждения, получим, что существует точка  $\xi \in [a, b]$  такая, что  $g^{(n+1)}(\xi) = 0$ . Вычисляя производную порядка  $(n + 1)$  функции  $g$ , получим:  $f^{(n+1)}(\xi) - K(n + 1)! = 0$ , т. е.  $K = f^{(n+1)}(\xi)/(n + 1)!$ , откуда вследствие (4.2) вытекает (4.1).  $\square$

Из только что доказанной теоремы вытекает простое, но полезное

**Следствие 4.1.** Пусть  $|f^{(n+1)}(x)| \leq M_{n+1} \forall x \in [a, b]$ . Тогда

$$|f(x) - L_n(x)| \leq \frac{M_{n+1}}{(n + 1)!} |\omega_{n+1}(x)| \quad \forall x \in [a, b]. \quad (4.3)$$

**4.1. Связь разделенных разностей с производными. Обобщенная формула Тейлора.** Используя представление (4.1) для погрешности интерполяционного полинома, установим связь между разделенными разностями и производными.

**Теорема 4.2.** Пусть функция  $f$  дифференцируема  $k$  раз на отрезке  $[\underline{x}, \bar{x}]$ , где  $\underline{x} = \min(x_0, x_1, \dots, x_k)$ ,  $\bar{x} = \max(x_0, x_1, \dots, x_k)$ . Тогда существует точка  $\xi \in [\underline{x}, \bar{x}]$  такая, что

$$f(x_0, x_1, \dots, x_k) = \frac{f^{(k)}(\xi)}{k!}. \quad (4.4)$$

**ДОКАЗАТЕЛЬСТВО.** Для упрощения записей будем считать, что узлами, по которым образуется разделенная разность, являются точки  $x_0, x_1, \dots, x_n, x$ , и построим интерполяционный полином  $P_{n+1}$  степени  $n + 1$  по значениям функции  $f$  в этих точках. Записывая значение полинома  $P_{n+1}$  в точке  $x$  в форме Ньютона, получим

$$P_{n+1}(x) = f(x_0) + (x - x_0)f(x_0, x_1) + \dots + \\ + (x - x_0)(x - x_1) \cdots (x - x_{n-1})f(x_0, x_1, \dots, x_n) +$$

$$\begin{aligned} &+ (x - x_0)(x - x_1) \cdots (x - x_n) f(x_0, x_1, \dots, x_n, x) = \\ &= P_n(x) + \omega_{n+1}(x) f(x_0, x_1, \dots, x_n, x) \end{aligned} \quad (4.5)$$

(см. (3.5)). С другой стороны, по построению  $P_{n+1}(x) = f(x)$ , следовательно,  $f(x) - P_n(x) = \omega_{n+1}(x) f(x_0, x_1, \dots, x_n, x)$ . Сравнивая это равенство с (4.1), получим утверждение теоремы.  $\square$

Выразим разделенную разность в каждом слагаемом равенства (4.5) через значение соответствующей производной. Тогда получим

$$\begin{aligned} f(x) = & f(x_0) + (x - x_0) f'(\xi_1) + (x - x_0)(x - x_1) \frac{f''(\xi_2)}{2!} + \dots + \\ & + (x - x_0)(x - x_1) \cdots (x - x_{n-1}) \frac{f^{(n)}(\xi_n)}{n!} + \omega_{n+1}(x) \frac{f^{(n+1)}(\xi_{n+1})}{(n+1)!}. \end{aligned} \quad (4.6)$$

Формулу (4.6) можно назвать обобщенной формулой Тейлора с остаточным членом в форме Лагранжа. В предельном случае, когда все узлы интерполирования стягиваются в узел  $x_0$ , формула принимает вид:

$$\begin{aligned} f(x) = & f(x_0) + (x - x_0) f'(x_0) + (x - x_0)^2 \frac{f''(x_0)}{2!} + \dots + \\ & + (x - x_0)^n \frac{f^{(n)}(x_0)}{n!} + (x - x_0)^{n+1} \frac{f^{(n+1)}(\xi)}{(n+1)!}, \end{aligned}$$

где  $\xi$  — некоторая точка, лежащая между  $x$  и  $x_0$ , т. е. (4.6) переходит в обычную формулу Тейлора с остаточным членом в форме Лагранжа.

**5. Оптимальный выбор узлов интерполирования. Многочлены Чебышева.** Из оценки (4.3) следует, что погрешность интерполирования зависит от расположения узлов.

Естественно попытаться решить следующую задачу: при заданном  $n$  расположить узлы  $x_0, x_1, \dots, x_n$  так, чтобы минимизировать величину

$$\max_{a \leq x \leq b} |(x - x_0)(x - x_1) \cdots (x - x_n)|. \quad (5.1)$$

При решении этой задачи будем использовать многочлены Чебышева, определяемые при помощи формул:

$$T_0(x) \equiv 1, \quad T_1(x) = x, \quad (5.2)$$

$$T_{k+1} = 2xT_k(x) - T_{k-1}(x), \quad k = 1, 2, \dots \quad (5.3)$$

Здесь  $k$  — степень многочлена.



Нетрудно видеть, что  $T_k(x) = 2^{k-1}x^k + \dots$ . Многочлены  $\tilde{T}_k(x) = 2^{1-k}T_k(x)$  со старшим коэффициентом, равным единице, называют нормированными многочленами Чебышева.

Построим явную формулу для полиномов Чебышева. Будем разыскивать значение  $T_k(x)$  в виде  $T_k(x) = \lambda^k$ . Используя это представление в рекуррентной формуле (5.3), получим  $\lambda^{k+1} = 2x\lambda^k - \lambda^{k-1}$ , откуда, предполагая, что  $\lambda \neq 0$ , приходим к квадратному уравнению  $\lambda^2 - 2x\lambda + 1 = 0$  для определения параметра  $\lambda$ . Корни этого уравнения:  $\lambda_{1,2} = x \pm \sqrt{x^2 - 1}$ , поэтому функции  $T_k^{(1)}(x) = (x + \sqrt{x^2 - 1})^k$ ,  $T_k^{(2)}(x) = (x - \sqrt{x^2 - 1})^k$ , а следовательно, и функции  $T_k(x) = c_1 T_k^{(1)}(x) + c_2 T_k^{(2)}(x)$ ,  $k = 0, 1, \dots$ , где  $c_1, c_2$  — произвольные постоянные, удовлетворяют рекуррентному соотношению (5.3). Выберем  $c_1, c_2$  так, чтобы были выполнены начальные условия (5.2):

$$c_1 + c_2 = 1,$$

$$(c_1 + c_2)x + (c_1 - c_2)\sqrt{x^2 - 1} = x.$$

Отсюда  $c_1 = c_2 = 1/2$ , т. е. полиномы

$$T_k(x) = \frac{1}{2} \left( x + \sqrt{x^2 - 1} \right)^k + \frac{1}{2} \left( x - \sqrt{x^2 - 1} \right)^k, \quad k = 0, 1, 2, \dots$$

удовлетворяют рекуррентному соотношению (5.3) и начальным условиям (5.2). При  $|x| \leq 1$  полиномам Чебышева можно придать более компактный вид. Положим в этом случае  $x = \cos \varphi$ . Тогда

$$T_k(x) = \frac{1}{2} (\cos \varphi + i \sin \varphi)^k + \frac{1}{2} (\cos \varphi - i \sin \varphi)^k,$$

откуда, используя формулы Муавра, получим  $T_k(x) = \cos k\varphi$ , или

$$T_k(x) = \cos k \arccos x. \quad (5.4)$$

Представление (5.4) позволяет найти все корни полинома Чебышева. В самом деле, из уравнения  $\cos k \arccos x = 0$  сразу получаем, что  $k \arccos x = \pi(2j + 1)/2$ ,  $j = 0, \pm 1, \pm 2, \dots$ , следовательно, точки

$$x_j = \cos \frac{\pi(2j + 1)}{2k}, \quad j = 0, 1, 2, \dots, k - 1,$$

есть корни полинома  $T_k(x)$ .

Найдем также точки экстремума полинома  $T_k(x)$  на отрезке  $[-1, 1]$ . Ясно, что — это те точки, в которых  $|T_k(x)| = 1$ . Они определяются соотношениями

$$\hat{x}_k = \cos \frac{j\pi}{k}, \quad j = 0, 1, \dots, k,$$

причем  $T_k(\hat{x}_j) = (-1)^j$ ,  $j = 0, 1, \dots, k$ .

Полиномы Чебышева являются полиномами, наименее уклоняющимися от нуля на отрезке  $[-1, 1]$ . А именно, справедлива

**Теорема 5.1.** *Для любого полинома  $Q_k(x)$  степени  $k$  со старшим коэффициентом, равным единице, справедливо неравенство*

$$\max_{-1 \leq x \leq 1} |Q_k(x)| > \max_{-1 \leq x \leq 1} |\tilde{T}_k(x)|,$$

где  $\tilde{T}_k(x)$  — нормированный полином Чебышева степени  $k$ .

**ДОКАЗАТЕЛЬСТВО.** Предположим противное, и пусть  $Q_k(x)$  — полином степени  $k$  со старшим коэффициентом, равным единице, для которого

$$\max_{-1 \leq x \leq 1} |Q_k(x)| \leq \max_{-1 \leq x \leq 1} |\tilde{T}_k(x)|.$$

Ясно, что полином  $R(x) = \tilde{T}_k(x) - Q_k(x)$  имеет степень не выше  $k - 1$ . Проанализируем значения полинома  $R(x)$  в точках экстремума полинома  $\tilde{T}_k(x)$ . По определению полинома  $\tilde{T}_k(x)$  имеем  $\tilde{T}_k(\hat{x}_j) = (-1)^j 2^{1-k}$ , а вследствие выдвинутого нами предположения  $|Q_k(\hat{x}_j)| \leq 2^{1-k}$ , поэтому  $R(\hat{x}_0) \geq 0$ ,  $R(\hat{x}_1) \leq 0$ ,  $R(\hat{x}_2) \geq 0, \dots$  Таким образом, на каждом отрезке  $[\hat{x}_l, \hat{x}_{l+1}]$ ,  $l = 0, 1, \dots, k - 1$ , полином  $R(x)$  имеет по крайней мере по одному корню, но это противоречит тому, что  $R(x)$  — полином степени не выше, чем  $k - 1$ .  $\square$

Используя полученный результат, нетрудно построить полином, наименее уклоняющийся от нуля на произвольном отрезке  $[a, b]$ . С этой целью выполним линейную замену переменной, переводящую отрезок  $-1 \leq t \leq 1$  в отрезок  $a \leq x \leq b$ :

$$x = \frac{a+b}{2} - \frac{a-b}{2}t.$$

При указанной замене полином  $T_k(t)$  переходит в полином

$$T_k\left(\frac{2x - (b+a)}{b-a}\right).$$

Коэффициент при  $x^k$  у этого полинома равен  $2^{2k-1}/(b-a)^k$ , а корнями являются точки

$$x_j = \frac{a+b}{2} + \frac{b-a}{2} \cos \frac{\pi(2j+1)}{2k}, \quad j = 0, 1, \dots, k-1.$$

Понятно, что полиномом, наименее уклоняющимся от нуля на отрезке  $[a, b]$ , среди всех полиномов степени  $k$  со старшим коэффициентом, равным единице, является полином

$$\frac{(b-a)^k}{2^{2k-1}} T_k \left( \frac{2x - (b+a)}{b-a} \right).$$

Таким образом, оптимальными узлами интерполирования, минимизирующими величину (5.1), являются точки

$$x_j = \frac{a+b}{2} + \frac{b-a}{2} \cos \frac{(2j+1)\pi}{2(n+1)}, \quad j = 0, 1, \dots, n. \quad (5.5)$$

Решение задачи об оптимальном расположении узлов интерполирования иллюстрирует рис. 1. Сплошной линией здесь изображен график функции  $\omega_9(x)$  на отрезке  $0 \leq x \leq 1$  при равномерном распределении узлов интерполирования, график функции  $\omega_9(x)$  при

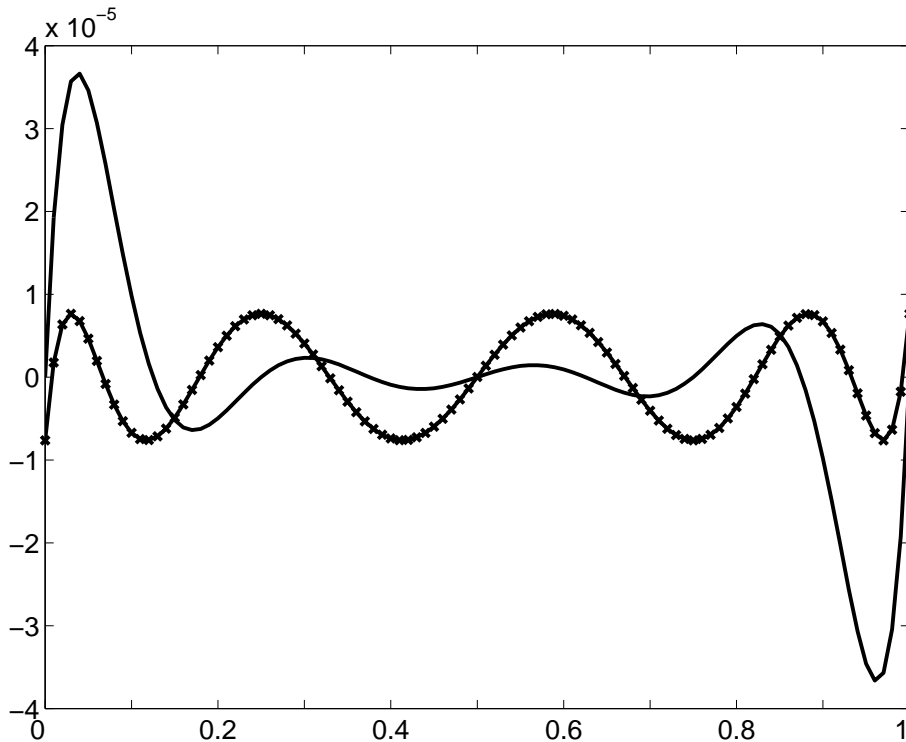


Рис. 1. Погрешность интерполирования.

оптимальном распределении узлов в соответствии с (5.5) помечен значком «×». При равномерном распределении узлов интерполирования погрешность сильно возрастает вблизи границ отрезка.

**6. Интерполирование с кратными узлами.** Будем предполагать теперь, что наряду со значениями функции  $f$  в узлах интерполирования заданы значения производных этой функции вплоть до определенного порядка. Точнее говоря, требуется найти такой полином  $P_m(x)$  степени  $m$ , что

$$P_m^{(j)}(x_i) = f^{(j)}(x_i), \quad j = 0, 1, \dots, m_i - 1, \quad i = 0, 1, \dots, n. \quad (6.1)$$

Число  $m_i \geq 1$  называется кратностью узла  $x_i$ . Понятно, что если  $m_i = 1$  при  $i = 0, 1, \dots, n$ , то есть все узлы простые, то получаем обычную задачу интерполирования.

**Теорема 6.1.** Пусть

$$m + 1 = m_0 + m_1 + \dots + m_n. \quad (6.2)$$

Тогда полином  $P_m(x)$  степени  $m$ , удовлетворяющий условиям (6.1), при любых значениях  $f^{(j)}(x_i)$ ,  $j = 0, 1, \dots, m_i - 1$ ,  $i = 0, 1, \dots, n$ , существует и определяется однозначно.

**ДОКАЗАТЕЛЬСТВО.** Условия (6.1) представляют собой систему линейных алгебраических уравнений относительно коэффициентов полинома  $P_m(x)$ . Условие (6.2) показывает, что число уравнений и число неизвестных в этой системе одно и то же, поэтому для доказательства теоремы достаточно убедиться, что соответствующая однородная система уравнений

$$P_m^{(j)}(x_i) = 0, \quad j = 0, 1, \dots, m_i - 1, \quad i = 0, 1, \dots, n \quad (6.3)$$

имеет только тривиальное решение. Равенства же (6.3) означают, что точки  $x_i$  — корни полинома  $P_m(x)$  с кратностями  $m_i$ , то есть полином  $P_m(x)$  степени  $m$  имеет  $m + 1$  корень, что может быть выполнено лишь в том случае, когда все его коэффициенты равны нулю.  $\square$

Как и в случае интерполирования с простыми узлами, можно выписать явный вид полинома  $P_m(x)$  либо в форме, аналогичной интерполяционному полиному Ньютона, либо в форме, аналогичной интерполяционному полиному Лагранжа. В последнем случае приходят к так называемой интерполяционной формуле Эрмита:

$$H_m(x) = \sum_{i=0}^n \sum_{j=0}^{m_i-1} f^{(j)}(x_i) \varphi_{ij}(x). \quad (6.4)$$

Здесь  $\varphi_{ij}(x)$  — базисные функции Эрмита, а именно полиномы степени  $m$ , удовлетворяющие условиям:

$$\varphi_{ij}^{(k)}(x_l) = \begin{cases} 0, & (l-i)^2 + (j-k)^2 \neq 0; \\ 1, & (l-i)^2 + (j-k)^2 = 0. \end{cases} \quad (6.5)$$

К сожалению, явный вид полиномов  $\varphi_{ij}(x)$  в произвольном случае оказывается весьма громоздким. Ограничимся поэтому рассмотрением следующего примера.

В точках  $x_0 < x_1 < x_2$  заданы значения  $f(x_0), f(x_1), f'(x_1), f(x_2)$ . Требуется построить полином третьей степени  $P_3(x)$  такой, что

$$P_3(x_0) = f(x_0), P_3(x_1) = f(x_1), P_3'(x_1) = f'(x_1), P_3(x_2) = f(x_2). \quad (6.6)$$

Будем искать полином  $P_3(x)$  в форме Эрмита:

$$P_3(x) \equiv H_3(x) = f(x_0)\varphi_{00}(x) + f(x_1)\varphi_{10}(x) + f'(x_1)\varphi_{11}(x) + f(x_2)\varphi_{20}(x).$$

При этом требуется построить базисные функции Эрмита, т. е. полиномы третьей степени, удовлетворяющие условиям:

$$\varphi_{00}(x_0) = 1, \varphi_{00}(x_1) = 0, \varphi'_{00}(x_1) = 0, \varphi_{00}(x_2) = 0, \quad (6.7)$$

$$\varphi_{10}(x_0) = 0, \varphi_{10}(x_1) = 1, \varphi'_{10}(x_1) = 0, \varphi_{10}(x_2) = 0, \quad (6.8)$$

$$\varphi_{11}(x_0) = 0, \varphi_{11}(x_1) = 0, \varphi'_{11}(x_1) = 1, \varphi_{11}(x_2) = 0, \quad (6.9)$$

$$\varphi_{20}(x_0) = 0, \varphi_{20}(x_1) = 0, \varphi'_{20}(x_1) = 0, \varphi_{20}(x_2) = 1. \quad (6.10)$$

Условия (6.7) показывают, что  $x_1, x_2$  — корни полинома  $\varphi_{00}$ , причем  $x_1$  есть корень кратности два, поэтому  $\varphi_{00}(x) = A(x - x_1)^2(x - x_2)$ , где  $A$  — постоянная, которую определим, используя первое из условий (6.7). Таким образом,

$$\varphi_{00}(x) = \frac{(x - x_1)^2(x - x_2)}{(x_0 - x_1)^2(x_0 - x_2)}.$$

Аналогично, используя условия (6.10), (6.9), получаем

$$\varphi_{20}(x) = \frac{(x - x_0)(x - x_1)^2}{(x_2 - x_0)(x_2 - x_1)^2},$$

$$\varphi_{11}(x) = \frac{(x - x_0)(x - x_1)(x - x_2)}{(x_1 - x_2)(x_1 - x_0)}.$$

Условия (6.8) определяют только два корня полинома  $\varphi_{10}(x)$ , поэтому  $\varphi_{10}(x) = (x - x_0)(x - x_2)(\alpha x + \beta)$ . Постоянные  $\alpha, \beta$  найдем из условий

$$\varphi_{10}(x_1) = (x_1 - x_0)(x_1 - x_2)(\alpha x_1 + \beta) = 1,$$

$$\varphi'_{10}(x_1) = (x_1 - x_0)(x_1 - x_2)\alpha + (\alpha x_1 + \beta)(2x_1 - x_0 - x_2) = 0.$$

Решая эту систему уравнений, получим

$$\alpha = \frac{2x_1 - x_0 - x_2}{(x_1 - x_0)^2(x_1 - x_2)^2},$$

$$\beta = \frac{1}{(x_1 - x_0)(x_1 - x_2)} \left( 1 + \frac{(2x_1 - x_0 - x_2)x_1}{(x_1 - x_0)(x_1 - x_2)} \right).$$

Таким образом, все базисные функции Эрмита, а следовательно, и полином  $H_3(x)$ , построены.

Оценку погрешности интерполяционного полинома с кратными узлами дает

**Теорема 6.2.** Пусть  $[a, b]$  — отрезок вещественной оси, содержащий все узлы интерполирования  $x_0, x_1, \dots, x_n$ , функция  $f$  непрерывно дифференцируема  $m + 1$  раз на отрезке  $[a, b]$ . Пусть, далее, точка  $x$  принадлежит  $[a, b]$  и не совпадает ни с одним из узлов интерполирования. Тогда существует точка  $\xi \in [a, b]$  такая, что

$$f(x) = P_m(x) + \frac{f^{(m+1)}(\xi)}{(m+1)!} \omega(x), \quad (6.11)$$

где  $P_m$  — полином степени  $m = m_0 + m_1 + \dots + m_n - 1$ , определяемый условиями (6.1),  $\omega(x) = (x - x_0)^{m_0} \dots (x - x_n)^{m_n}$ .

Доказательство этой теоремы аналогично доказательству теоремы 4.1.

## § 2. Среднеквадратичное приближение функций

**1. Элемент наилучшего среднеквадратичного приближения.** Говорят, что функции  $f$  и  $\varphi$  близки в среднеквадратичном смысле на отрезке  $[a, b]$ , если мал интеграл

$$\int_a^b p(x)(f(x) - \varphi(x))^2 dx.$$

Здесь  $p$  — заданная положительная на отрезке  $[a, b]$  функция, называемая весом.

Пусть  $\varphi_1(x), \varphi_2(x), \dots, \varphi_n(x)$  — заданные на отрезке  $[a, b]$  функции. Функция  $\varphi(x) = c_1\varphi_1(x) + c_2\varphi_2(x) + \dots + c_n\varphi_n(x)$  называется элементом наилучшего среднеквадратичного приближения к функции  $f(x)$ , если числа  $c_1, \dots, c_n$  доставляют минимальное значение интегралу

$$\int_a^b p(x) \left( f(x) - \sum_{k=1}^n c_k \varphi_k(x) \right)^2 dx.$$

При исследовании существования и единственности элемента наилучшего приближения и при описании способов его построения будет полезно ввести некоторые новые понятия и обозначения.

Величину

$$(f, g) = \int_a^b p(x)f(x)g(x)dx$$

называют скалярным произведением функций  $f, g$ . Скалярное произведение функций обладает свойствами, аналогичными свойствам скалярного произведения векторов:

- 1)  $(f, g) = (g, f)$ ,
- 2)  $(f, f) > 0$  для любой функции  $f$ , не равной тождественно нулю.
- 3)  $(\alpha_1 f_1 + \alpha_2 f_2, g) = \alpha_1 (f_1, g) + \alpha_2 (f_2, g)$  для любых функций  $f_1, f_2, g$  и чисел  $\alpha_1, \alpha_2$ .

Число  $\|f\| = (f, f)^{1/2}$  называется нормой функции  $f$ . Норма обладает свойствами, аналогичными свойствам длины вектора.

Напомним определение линейной независимости функций. Функции  $\varphi_1(x), \dots, \varphi_n(x)$  называются линейно независимыми на отрезке  $[a, b]$ , если равенство

$$c_1\varphi_1(x) + c_2\varphi_2(x) + \dots + c_n\varphi_n(x) = 0 \quad \forall x \in [a, b]$$

выполнено лишь при условии, что  $c_1 = c_2 = \dots = c_n = 0$ .

Введем в рассмотрение вещественную функцию

$$F(c) = \left\| f - \sum_{i=1}^n c_i \varphi_i \right\|^2$$

$n$  независимых переменных  $c_1, \dots, c_n$ . Задача построения элемента наилучшего среднеквадратичного приближения эквивалентна задаче минимизации функции  $F$ .

Выпишем необходимые условия минимума функции  $F$ . Предварительно, используя свойства симметрии и линейности скалярного произведения, преобразуем ее к виду:

$$F(c) = (f, f) + \sum_{i,j=1}^n c_i c_j (\varphi_i, \varphi_j) - 2 \sum_{i=1}^n c_i (f, \varphi_i).$$

Вычисляя первые производные, получим

$$\frac{\partial F(c)}{\partial c_i} = 2 \sum_{j=1}^n c_j (\varphi_j, \varphi_i) - 2(f, \varphi_i), \quad i = 1, 2, \dots, n.$$

Таким образом, система уравнений для отыскания компонент вектора  $c$ , доставляющего минимальное значение функции  $F$ , имеет вид

$$\sum_{j=1}^n c_j (\varphi_j, \varphi_i) = (f, \varphi_i), \quad i = 1, 2, \dots, n. \quad (1.1)$$

Матрица этой системы линейных уравнений  $A = \{(\varphi_i, \varphi_j)\}_{i,j=1}^n$  называется матрицей Грама системы функций  $\varphi_1, \dots, \varphi_n$ . Матрица  $A$ , очевидно, симметрична.

**Теорема 1.1.** *Для того, чтобы матрица Грама была положительно определена, необходимо и достаточно, чтобы функции  $\varphi_1, \dots, \varphi_n$  были линейно независимы.*

**ДОКАЗАТЕЛЬСТВО.** Выполненные выше преобразования функции  $F$  показывают, что  $(Ac, c) = \left\| \sum_{i=1}^n c_i \varphi_i \right\|^2 \geq 0$ . Если функции  $\varphi_1, \dots, \varphi_n$  линейно независимы, то равенство в последнем неравенстве достигается лишь при  $c_1 = c_2 = \dots = c_n = 0$ , т. е. матрица  $A$  положительно определена. Обратное утверждение теоремы доказывается аналогично.  $\square$

Из только что доказанной теоремы, очевидно, вытекает, что если функции  $\varphi_1, \dots, \varphi_n$  линейно независимы, то элемент наилучшего приближения для любой функции  $f$  существует и определяется единственным образом.

Наиболее просто элемент наилучшего приближения строится в том случае, когда матрица  $A$  диагональна, т. е., когда  $(\varphi_i, \varphi_j) = 0$  при  $i \neq j$ . Система функций  $\varphi_1, \dots, \varphi_n$  называется в этом случае ортогональной, а элемент наилучшего приближения принимает вид

$$\varphi(x) = \sum_{i=1}^n \frac{(f, \varphi_i)}{(\varphi_i, \varphi_i)} \varphi_i(x).$$

Чаще всего на практике применяются ортогональные системы функций, состоящие из тригонометрических функций или из полиномов.

Например, нетрудно проверить, что

$$\int_0^{2\pi} \sin kx \, dx = 0, \quad \int_0^{2\pi} \cos kx \, dx = 0, \quad k = 1, 2, \dots$$



$$\int_0^{2\pi} \sin kx \cos lx \, dx = 0, \quad k, l = 1, 2, \dots$$

$$\int_0^{2\pi} \sin kx \sin lx \, dx = \begin{cases} 0, & k \neq l; \\ \pi, & k = l \end{cases}; \quad \int_0^{2\pi} \cos kx \cos lx \, dx = \begin{cases} 0, & k \neq l; \\ \pi, & k = l. \end{cases}$$

Поэтому функции  $1, \sin x, \cos x, \sin 2x, \cos 2x, \dots, \sin nx, \cos nx$  при любом  $n \geq 1$  образуют ортогональную систему на отрезке  $[0, 2\pi]$ .

Элемент наилучшего среднеквадратичного приближения по этой системе функций имеет вид

$$\varphi(x) = a_0 + \sum_{k=1}^n (a_k \cos kx + b_k \sin kx),$$

где

$$a_0 = \frac{1}{2\pi} \int_0^{2\pi} f(x) \, dx, \quad a_k = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos kx \, dx, \quad b_k = \frac{1}{\pi} \int_0^{2\pi} f(x) \sin kx \, dx,$$

$k = 1, 2, \dots, n$ .

Функцию  $\varphi(x)$  называют в этом случае отрезком ряда Фурье для функции  $f(x)$ , числа  $a_k, b_k$  — коэффициентами Фурье функции  $f(x)$ .

**2. Ортогональные полиномы.** Подробнее рассмотрим ортогональные системы функций, образованные полиномами  $P_0(x), P_1(x), \dots, P_n(x)$ .

**Теорема 2.1.** *Для любого заданного отрезка  $[a, b]$  и веса  $p(x) > 0$  система ортогональных полиномов существует.*

**ДОКАЗАТЕЛЬСТВО.** Построим искомую систему полиномов  $P_0(x), P_1(x), P_2(x), \dots, P_n(x)$ , применяя метод ортогонализации Грама — Шмидта к системе степеней независимой переменной:  $1, x, x^2, \dots, x^n$ . Положим  $P_0(x) \equiv 1$ . Полином  $P_1(x)$  будем разыскивать в виде  $P_1(x) = x - \alpha_{10}P_0(x)$ . Постоянную  $\alpha_{10}$  определим из условия ортогональности  $(P_1, P_0) = 0$ . Получим:  $\alpha_{10} = (x, P_0)/(P_0, P_0)$ . Вообще, если полиномы  $P_0(x), P_1(x), \dots, P_k(x)$  уже построены, то полином  $P_{k+1}(x)$  будем разыскивать в виде

$$P_{k+1}(x) = x^{k+1} - \sum_{i=1}^k \alpha_{k+1,i} P_i(x),$$

определяя постоянные  $\alpha_{k+1,i}$  из условий ортогональности  $(P_{k+1}, P_i) = 0$ ,  $i = 1, 2, \dots, k$ . В результате получим  $\alpha_{k+1,i} = (x^{k+1}, P_i) / (P_i, P_i)$ .  $\square$

Отметим некоторые свойства ортогональных полиномов, непосредственно вытекающие из определения ортогональности.

1. Произвольный полином  $R_k(x)$  степени  $k$  однозначно представим в виде линейной комбинации ортогональных полиномов:

$$R_k(x) = \sum_{i=0}^k \alpha_i P_i(x),$$

где  $\alpha_i = (R_k, P_i) / (P_i, P_i)$ ,  $i = 0, 1, \dots, n$ .

2. Ортогональный полином степени  $k$  ортогонален любому полиному  $q_i$  меньшей степени:  $(P_k, q_i) = 0$ ,  $i = 0, 1, \dots, k-1$ <sup>1)</sup>.

3. Задание веса и отрезка ортогональности определяет систему ортогональных полиномов с точностью до множителя, то есть если  $Q_0(x), Q_1(x), \dots$  ортогональны, как и полиномы  $P_0(x), P_1(x), \dots$ , с весом  $p(x)$  на отрезке  $[a, b]$ , то  $Q_i(x) = \alpha_i P_i(x)$ ,  $i = 0, 1, \dots$ , где  $\alpha_i$  — ненулевые постоянные.

4. Ортогональные полиномы удовлетворяют трехчленному рекуррентному соотношению:

$$(x - \beta_{k+1,k})P_k(x) = \beta_{k+1,k+1}P_{k+1}(x) + \beta_{k+1,k-1}P_{k-1}(x), \quad (2.1)$$

$k = 1, 2, \dots$

Докажем лишь последнее утверждение. Остальные доказываются проще. Представим полином  $xP_k(x)$  в виде разложения по ортогональным полиномам:

$$xP_k(x) = \beta_{k+1,k+1}P_{k+1}(x) + \beta_{k+1,k}P_k(x) + \beta_{k+1,k-1}P_{k-1}(x) + \dots + \beta_{k+1,0}P_0(x),$$

где  $\beta_{k+1,i} = (xP_k, P_i) / (P_i, P_i)$ ,  $i = 0, 1, \dots, k+1$ . Записывая более подробно числитель последнего выражения, получим

$$(xP_k, P_i) = \int_a^b p(x)xP_k(x)P_i(x) dx = (P_k, q_{i+1}),$$

где  $q_{i+1}(x)$  — полином степени  $i+1$ . Отсюда вследствие свойства 2) вытекает, что  $\beta_{k+1,i} = 0$  при  $i < k-1$ .  $\square$

<sup>1)</sup>Нетрудно убедиться, что это свойство можно принять за определение ортогонального полинома.

**Теорема 2.2.** Пусть  $P_n(x)$  — полином, ортогональный на отрезке  $[a, b]$ . Тогда все его корни вещественны, различны и принадлежат интервалу  $(a, b)$ .

**ДОКАЗАТЕЛЬСТВО.** Обозначим через  $x_1, x_2, \dots, x_k$  точки, принадлежащие  $(a, b)$ , в которых полином  $P_n(x)$  меняет знак. Предположение о том, что  $k < n$  приводит к противоречию. Действительно, в этом случае  $P_n(x) = (x - x_1)(x - x_2) \dots (x - x_k)Q_{n-k}(x)$ , где  $Q_{n-k}(x)$  — полином степени  $n - k$ . Ясно, что корнями полинома  $Q_{n-k}(x)$  по построению являются корни полинома  $P_n(x)$ , принадлежащие  $(a, b)$  и имеющие четную кратность, а также корни полинома  $P_n(x)$ , не принадлежащие  $(a, b)$ . Поэтому полином  $Q_{n-k}(x)$  не меняет знака на отрезке  $[a, b]$ . Положим  $q_k(x) = (x - x_1)(x - x_2) \dots (x - x_k)$ . Тогда

$$(P_n, q_k) = \int_a^b p(x)q_k^2(x)Q_{n-k}(x) dx \neq 0,$$

поскольку подынтегральное выражение не меняет знака на интервале  $(a, b)$  и лишь в конечном числе точек обращается в нуль. Но, с другой стороны,  $(P_n, q_k) = 0$ , так как полином  $P_n(x)$  ортогонален, а  $k < n$ .  $\square$

Корни ортогональных полиномов обладают весьма важным свойством перемежаемости. Сформулируем его без доказательства.

**Теорема 2.3.** Пусть  $P_n(x), P_{n+1}(x), n \geq 1$  — полиномы, ортогональные на отрезке  $[a, b]$ . Обозначим через  $x_1^{(n)} < x_2^{(n)} < \dots < x_n^{(n)}$ ,  $x_1^{(n+1)} < x_2^{(n+1)} < \dots < x_{n+1}^{(n+1)}$  их корни. Тогда

$$a < x_1^{(n+1)} < x_1^{(n)} < x_2^{(n+1)} < x_2^{(n)} < \dots < x_n^{(n)} < x_{n+1}^{(n+1)} < b.$$

Примеры ортогональных полиномов.

1. Полиномы Чебышева ортогональны на отрезке  $[-1, 1]$  с весом  $p(x) = 1/\sqrt{1-x^2}$ . Действительно, используя представление (5.4), получим

$$(T_k, T_l) = \int_{-1}^1 \frac{\cos(k \arccos x) \cos(l \arccos x)}{\sqrt{1-x^2}} dx.$$

Полагая  $x = \cos \varphi$ , нетрудно подсчитать, что

$$(T_k, T_l) = \int_0^\pi \cos k\varphi \cos l\varphi d\varphi = \frac{1}{2} \int_0^\pi (\cos(k+l)\varphi + \cos(k-l)\varphi) d\varphi = 0$$

при  $k \neq l$ .

2. Полиномы Лежандра. Построим полиномы, ортогональные с весом  $p(x) \equiv 1$  на отрезке  $[-1, 1]$ . Такие полиномы называют полиномами Лежандра. Обозначим искомым полином через  $L_n(x)$ ,  $n \geq 1$ . По определению

$$\int_{-1}^1 L_n(x) q_{n-1}(x) dx = 0 \quad (2.2)$$

для любого полинома  $q_{n-1}(x)$  степени  $n - 1$ .

Положим  $L_n(x) = Q_{2n}^{(n)}(x)$ , где  $Q_{2n}(x)$  — некоторый полином степени  $2n$ , и преобразуем интеграл в правой части равенства (2.2) при помощи формулы интегрирования по частям. Получим

$$\int_{-1}^1 Q_{2n}^{(n)}(x) q_{n-1}(x) dx = - \int_{-1}^1 Q_{2n}^{(n-1)}(x) q'_{n-1}(x) dx + Q_{2n}^{(n-1)}(x) q_{n-1}(x) \Big|_{-1}^1.$$

Продолжая аналогичные преобразования и учитывая, что  $q_{n-1}^{(n)}(x) \equiv 0$ , можем написать:

$$\begin{aligned} \int_{-1}^1 Q_{2n}^{(n)}(x) q_{n-1}(x) dx &= Q_{2n}^{(n-1)}(x) q_{n-1}(x) \Big|_{-1}^1 - Q_{2n}^{(n-2)}(x) q'_{n-1}(x) \Big|_{-1}^1 + \\ &+ Q_{2n}^{(n-3)}(x) q''_{n-1}(x) \Big|_{-1}^1 - \dots + (-1)^{n-1} Q_{2n}(x) q_{n-1}^{(n-1)}(x) \Big|_{-1}^1. \end{aligned}$$

Условие (2.2) будет выполнено, если положить, что

$$Q_{2n}(1) = Q'_{2n}(1) = \dots = Q_{2n}^{(n-1)}(1) = 0, \quad (2.3)$$

$$Q_{2n}(-1) = Q'_{2n}(-1) = \dots = Q_{2n}^{(n-1)}(-1) = 0. \quad (2.4)$$

Равенства (2.3), (2.4) означают, что точки  $-1, +1$  — корни кратности  $n$  полинома  $Q_{2n}(x)$ , т. е.  $Q_{2n}(x) = A(x-1)^n(x+1)^n$ , где  $A$  — произвольная постоянная. Таким образом, полином Лежандра представим в виде

$$L_n(x) = \frac{d^n}{dx^n} (1-x^2)^n. \quad (2.5)$$

Полученную формулу называют формулой Родрига.

### § 3. Приближенное вычисление интегралов

В этом параграфе будут рассмотрены методы приближенного вычисления определенных интегралов

$$\int_a^b f(x)dx.$$

Лишь в исключительных случаях удастся найти первообразную  $F(x)$  функции  $f(x)$  и представить интеграл в виде

$$\int_a^b f(x)dx = F(b) - F(a).$$

Как правило, для вычисления интегралов приходится применять приближенные методы, и чаще всего применяют так называемые формулы механических квадратур, т. е. разыскивают приближенное значение интеграла

$$\int_a^b f(x)dx \approx \sum_{k=1}^n c_k f(x_k).$$

Числа  $c_k$  называют весами, или коэффициентами квадратурной формулы,  $x_k$  — узлами квадратурной формулы.

Конкретные квадратурные формулы отличаются, таким образом, выбором коэффициентов и узлов.

**1. Интерполяционные квадратурные формулы.** При описании этого класса квадратурных формул интеграл будет удобно записывать в виде

$$I(f) = \int_a^b p(x)f(x)dx,$$

где  $p(x)$  — положительная на отрезке  $[a, b]$  функция, называемая весом.

Выберем на отрезке  $[a, b]$  точки  $a \leq x_1 < x_2 \dots < x_n \leq b$ . Построим полином  $L_{n-1}(x)$  степени  $n - 1$ , интерполирующий функцию  $f(x)$  по указанным узлам, и примем за приближенное значение интеграла

величину

$$S(f) = \int_a^b p(x)L_{n-1}(x)dx.$$

В дальнейшем будем пользоваться также следующим обозначением для погрешности квадратурной формулы  $R(f) = I(f) - S(f)$ . Величину  $R(f)$  будем называть также остаточным членом квадратурной формулы. Записывая интерполяционный полином  $L_{n-1}(x)$  в форме

Лагранжа (см. (2.3), § 1):  $L_{n-1}(x) = \sum_{k=1}^n f(x_k)\varphi_k(x)$ , получим

$$S(f) = \sum_{k=1}^n c_k f(x_k), \quad (1.1)$$

где

$$c_k = \int_a^b p(x)\varphi_k(x)dx, \quad k = 1, 2, \dots, n. \quad (1.2)$$

Здесь  $\varphi_k(x)$  — базисные функции Лагранжа, определяемые формулами (2.2), § 1.

Квадратурная формула вида (1.1) с коэффициентами, вычисляемыми по формулам (1.2), называется интерполяционной квадратурной формулой.

Следующая теорема дает другое, эквивалентное описание интерполяционных квадратурных формул.

**Теорема 1.1.** *Для того, чтобы квадратурная формула вида (1.1) была интерполяционной, необходимо и достаточно, чтобы она была точной на любом полиноме степени  $n - 1$ , то есть  $R(P_{n-1}) = 0$  для любого полинома  $P_{n-1}(x)$  степени  $n - 1$ .*

**ДОКАЗАТЕЛЬСТВО.** Пусть формула (1.1) является интерполяционной. Тогда для любого полинома  $P_{n-1}(x)$  имеем  $R(P_{n-1}) = I(P_{n-1}) - S(P_{n-1}) = I(P_{n-1}) - I(L_{n-1})$ , где  $L_{n-1}(x)$  — полином степени  $n - 1$ , интерполирующий полином  $P_{n-1}$  по узлам  $x_1, x_2, \dots, x_n$ . Вследствие теоремы 1.1, § 1,  $L_{n-1}(x) \equiv P_{n-1}(x)$ , т. е.  $R(P_{n-1}) = 0$ . Докажем обратное утверждение. Поскольку базисная функция Лагранжа  $\varphi_k(x)$  при любом  $k = 1, 2, \dots, n$  есть полином степени  $n - 1$ , то  $I(\varphi_k) = S(\varphi_k)$ ,

но  $S(\varphi_k) = \sum_{i=1}^n c_i \varphi_k(x_i) = c_k$ .  $\square$

**2. Устойчивость квадратурных формул.** Важным свойством квадратурных формул является их устойчивость по отношению к ошибкам вычисления подынтегральной функции. При вычислении значений функции  $f$  неизбежно возникают ошибки, например, за счет округлений. В результате, вместо значений  $f(x_k)$  получаем значения  $\tilde{f}(x_k) = f(x_k) + \varepsilon_k$ , и, соответственно, вместо  $S(f)$  значения  $\tilde{S}(f) = S(f) + \sum_{k=1}^n c_k \varepsilon_k$ .

Квадратурная формула называется устойчивой, если существует постоянная  $C$  такая, что равномерно по числу узлов  $\sum_{k=1}^n |c_k| \leq C$ .

Если предположить, что равномерно по числу узлов квадратурной формулы выполняется оценка  $|\varepsilon_k| \leq \varepsilon$ , то для устойчивой квадратурной формулы  $|S(f) - \tilde{S}(f)| \leq C\varepsilon$ , т. е. погрешность вычисления  $S(f)$  определяется лишь погрешностью вычисления значений подынтегральной функции.

Справедлив следующий простой, но практически очень важный, достаточный признак устойчивости интерполяционных квадратурных формул.

**Теорема 2.1.** *Если при любом  $n \geq 1$  все коэффициенты интерполяционной квадратурной формулы неотрицательны, то она устойчива.*

**ДОКАЗАТЕЛЬСТВО.** По определению при любом  $n \geq 1$  интерполяционная квадратурная формула точна на полиноме нулевой степени, следовательно, учитывая неотрицательность коэффициентов  $c_k$ , получим

$$\int_a^b p(x) dx = \sum_{k=1}^n c_k = \sum_{k=1}^n |c_k|. \quad \square$$

**3. Квадратурные формулы Ньютона — Котеса.** Так называют интерполяционные квадратурные формулы при  $p(x) \equiv 1$  и равномерно расположенных на отрезке  $[a, b]$  узлах:  $x_i = x_1 + (i-1)h$ ,  $i = 2, \dots, n$ ,  $h = \text{const}$ . Выполняя в формуле (1.2) замену переменной  $x = x_1 + th$ , нетрудно получить, что

$$c_k = (b-a)\hat{c}_k, \quad k = 1, 2, \dots, n,$$

где числа  $\hat{c}_k$  зависят только от  $k$  и  $n$ .

Приведем примеры простейших квадратурных формул Ньютона — Котеса.

1. Квадратурная формула левых прямоугольников. В этом случае  $n = 1$ ,  $x_1 = a$ ,

$$\int_a^b f(x) dx \approx (b - a)f(a). \quad (3.1)$$

2. Квадратурная формула правых прямоугольников,  $n = 1$ ,  $x_1 = b$ ,

$$\int_a^b f(x) dx \approx (b - a)f(b). \quad (3.2)$$

3. Квадратурная формула центральных прямоугольников,  $n = 1$ ,  $x_1 = c = (b + a)/2$ ,

$$\int_a^b f(x) dx \approx (b - a)f(c). \quad (3.3)$$

4. Квадратурная формула трапеций. В этом случае  $n = 2$ ,  $x_1 = a$ ,  $x_2 = b$ ,  $h = b - a$ ,  $\varphi_1(x) = (b - x)/(b - a)$ ,  $\varphi_2(x) = (a - x)/(a - b)$ ,  $c_1 = c_2 = (b - a)/2$ ,

$$\int_a^b f(x) dx \approx \frac{(b - a)}{2}(f(a) + f(b)). \quad (3.4)$$

5. Квадратурная формула Симпсона (парабол). В этом случае  $n = 3$ ,  $x_1 = a$ ,  $x_2 = c = (a + b)/2$ ,  $x_3 = b$ ,  $h = (b - a)/2$ ,

$$\varphi_1(x) = \frac{(c - x)(b - x)}{(c - a)(b - a)}, \quad \varphi_2(x) = \frac{(a - x)(b - x)}{(a - c)(b - c)}, \quad \varphi_3(x) = \frac{(a - x)(c - x)}{(a - b)(c - b)},$$

$$c_1 = c_3 = (b - a)/6, \quad c_2 = 2(b - a)/3,$$

$$\int_a^b f(x) dx \approx \frac{(b - a)}{6}(f(a) + 4f(c) + f(b)). \quad (3.5)$$

Формулы (3.1)–(3.5) имеют простую геометрическую интерпретацию. Применяя формулу (3.1), мы заменяем площадь криволинейной трапеции, т. е.

$$\int_a^b f(x) dx,$$



площадью прямоугольника с основанием  $b-a$  и высотой, равной  $f(a)$ , при использовании формулы (3.2) — площадью прямоугольника с тем же основанием и высотой, равной  $f(b)$ . При использовании формулы (3.3) высота прямоугольника равна  $f(c)$ . В случае формулы (3.4) площадь криволинейной трапеции заменяется площадью прямолинейной трапеции, ограниченной прямой, проходящей через точки  $(a, f(a))$ ,  $(b, f(b))$ , и, наконец, при использовании формулы (3.5) вычисляется площадь криволинейной трапеции, ограниченной параболой, проходящей через три точки:  $(a, f(a))$ ,  $(c, f(c))$ ,  $(b, f(b))$ .

Геометрически очевидно, что при достаточно плавно меняющейся на отрезке  $[a, b]$  функции  $f(x)$  наиболее точной должна быть формула парабол.

Формулы Ньютона — Котеса при  $n > 4$  редко используются на практике. Дело в том, что уже при  $n = 8$  среди коэффициентов формулы встречаются отрицательные, и сумма модулей коэффициентов неограниченно возрастает с увеличением  $n$ , т. е. формулы Ньютона — Котеса неустойчивы.

**4. Оценки точности простейших квадратурных формул Ньютона — Котеса.** В этом пункте, предполагая, что функция  $f$  дифференцируема на отрезке  $[a, b]$  достаточное число раз, получим представления для погрешностей квадратурных формул, построенных в предыдущем пункте.

1. Квадратурные формулы правых и левых прямоугольников. Для формулы левых прямоугольников по построению

$$R(f) = \int_a^b (f(x) - f(a)) dx.$$

Применяя формулу конечных приращений Лагранжа, получим

$$f(x) - f(a) = f'(\zeta(x))(x - a),$$

Нетрудно видеть, что функция

$$f'(\zeta(x)) = \frac{f(x) - f(a)}{x - a}$$

непрерывна на отрезке  $[a, b]$ . При  $x \neq a$  это следует из непрерывности функции  $f$ . Особенность при  $x = a$  устранима, поскольку

$$\lim_{x \rightarrow a} f'(\zeta(x)) = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a} = f'(a).$$

Таким образом, получаем, что

$$R(f) = \int_a^b f'(\zeta(x))(x - a) dx.$$

Функция  $x - a$  не меняет знака на отрезке  $[a, b]$ . Поэтому, применяя теорему о среднем, приходим к представлению остаточного члена для квадратурной формулы левых прямоугольников

$$R(f) = f'(\xi) \int_a^b (x - a) dx = f'(\xi) \frac{(b - a)^2}{2}, \quad \xi \in (a, b). \quad (4.1)$$

Совершенно аналогичные рассуждения приводят к представлению остаточного члена квадратурной формулы правых прямоугольников

$$R(f) = f'(\xi) \int_a^b (x - b) dx = -f'(\xi) \frac{(b - a)^2}{2}, \quad \xi \in (a, b). \quad (4.2)$$

2. Квадратурная формула трапеций. В этом случае

$$R(f) = \int_a^b (f(x) - L_1(x)) dx,$$

где  $L_1(x)$  — полином первой степени, удовлетворяющий условиям:  $L_1(a) = f(a)$ ,  $L_1(b) = f(b)$ . Используя формулу (4.1), § 1, для остаточного члена интерполяционного полинома, будем иметь

$$R(f) = \int_a^b \frac{f''(\zeta(x))}{2} (x - a)(x - b) dx,$$

причем функция  $(x - a)(x - b)$  не меняет знака на отрезке интегрирования, а функция  $f''(\zeta(x))$ , как нетрудно проверить при помощи правила Лопиталья, непрерывна на этом отрезке. Вновь применяя теорему о среднем, получим

$$R(f) = \frac{f''(\xi)}{2} \int_a^b (x - a)(x - b) dx = -\frac{f''(\xi)}{12} (b - a)^3, \quad \xi \in (a, b). \quad (4.3)$$

3. Квадратурная формула центральных прямоугольников. Рассуждая сначала, как и в случае квадратурной формулы левых прямоугольников, получим

$$R(f) = \int_a^b f'(\zeta(x))(x - c) dx,$$

но функция  $x - c$  меняет знак в точке  $c$ . Поэтому применить теорему о среднем здесь не удастся. В связи с этим введем в рассмотрение полином первой степени  $L_1(x) = f(c) + \alpha(x - c)$ . При любом значении постоянной  $\alpha$

$$\int_a^b L_1(x) dx = f(c)(b - a),$$

следовательно,

$$R(f) = \int_a^b (f(x) - L_1(x)) dx.$$

Выберем теперь  $\alpha$  так, чтобы  $L_1'(c) = f'(c)$ . Ясно, что для этого достаточно положить  $\alpha = f'(c)$ . При таком  $\alpha$  точка  $c$  оказывается узлом интерполирования кратности два, следовательно, применяя формулу (6.11), § 1, можно написать

$$f(x) - L_1(x) = \frac{f''(\zeta(x))}{2}(x - c)^2,$$

причем функция  $f''(\zeta(x))$  непрерывна по  $x$  на  $[a, b]$ . Теперь теорема о среднем применима, поэтому

$$R(f) = \frac{f''(\xi)}{2} \int_a^b (x - c)^2 dx = \frac{f''(\xi)}{24}(b - a)^3, \quad \xi \in (a, b). \quad (4.4)$$

4. Формула Симпсона. В этом случае

$$R(f) = \int_a^b (f(x) - L_2(x)) dx,$$

где  $L_2(x)$  — полином второй степени, удовлетворяющий условиям

$$L_2(a) = f(a), \quad L_2(c) = f(c), \quad L_2(b) = f(b).$$

На основании формулы (4.1), § 1, имеем

$$R(f) = \int_a^b \frac{f'''(\zeta(x))}{3!} (x-a)(x-c)(x-b) dx.$$

Функция  $f'''(\zeta(x))$  непрерывна по  $x$ , функция  $(x-a)(x-c)(x-b)$  меняет знак в точке  $c$ . Вновь для того, чтобы получить возможность применить теорему о среднем, повысим порядок интерполяционного полинома, а именно, введем в рассмотрение полином

$$L_3(x) = L_2(x) + \alpha(x-a)(x-c)(x-b).$$

Нетрудно проверить, что функция  $(x-a)(x-c)(x-b)$  нечетна относительно точки  $c$ , следовательно,

$$\int_a^b (x-a)(x-c)(x-b) dx = 0,$$

поэтому

$$R(f) = \int_a^b (f(x) - L_3(x)) dx.$$

При любом значении  $\alpha$  имеем  $L_3(a) = f(a)$ ,  $L_3(c) = f(c)$ ,  $L_3(b) = f(b)$ . Заметим теперь, что  $((x-a)(x-c)(x-b))'|_{x=c} \neq 0$ , следовательно, постоянную  $\alpha$  можно выбрать так, чтобы  $L_3'(c) = f'(c)$ . Но тогда на основании формулы (6.11), § 1, получим

$$\begin{aligned} R(f) &= \int_a^b \frac{f^{IV}(\zeta(x))}{4!} (x-a)(x-c)^2(x-b) dx = \\ &= -\frac{f^{IV}(\xi)}{2880} (b-a)^5, \quad \xi \in (a, b). \end{aligned} \quad (4.5)$$

**5. Составные квадратурные формулы.** На практике формулы (3.1)–(3.5) чаще всего используют как составные, т. е. предварительно отрезок  $[a, b]$  разбивают на частичные отрезки достаточно малой длины, затем на каждом частичном отрезке применяют ту или иную квадратурную формулу. Введем необходимые обозначения. Положим  $a = x_0 < x_1 < \dots < x_N = b$ ,  $h_i = x_i - x_{i-1}$ ,  $i = 1, 2, \dots, N$ ,  $h = \max_{1 \leq i \leq N} h_i$ ,  $x_{i-1/2} = x_i - h_i/2$ .

1. Формула левых прямоугольников:

$$\int_a^b f(x) dx = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x) dx \approx \sum_{i=1}^N f(x_{i-1})h_i. \quad (5.1)$$

Используя (4.1), получим

$$R(f) = \sum_{i=1}^N \left( \int_{x_{i-1}}^{x_i} f(x) dx - f(x_{i-1})h_i \right) = \frac{1}{2} \sum_{i=1}^N f'(\xi_i)h_i^2,$$

$\xi_i \in (x_{i-1}, x_i)$ ,  $i = 1, \dots, N$ , откуда

$$|R(f)| \leq \frac{M_1}{2} h \sum_{i=1}^N h_i = \frac{M_1}{2} h(b-a), \quad (5.2)$$

где  $M_1 = \max_{a \leq x \leq b} |f'(x)|$ .

2. Формула правых прямоугольников:

$$\int_a^b f(x) dx \approx \sum_{i=1}^N f(x_i)h_i. \quad (5.3)$$

Используя (4.2), получим

$$R(f) = -\frac{1}{2} \sum_{i=1}^N f'(\xi_i)h_i^2,$$

$\xi_i \in (x_{i-1}, x_i)$ ,  $i = 1, \dots, N$ , откуда

$$|R(f)| \leq \frac{M_1}{2} h(b-a). \quad (5.4)$$

3. Формула центральных прямоугольников:

$$\int_a^b f(x) dx \approx \sum_{i=1}^N f(x_{i-1/2})h_i. \quad (5.5)$$

Используя (4.4), получим

$$R(f) = \frac{1}{24} \sum_{i=1}^N f''(\xi_i)h_i^3,$$

$\xi_i \in (x_{i-1}, x_i)$ ,  $i = 1, \dots, N$ , откуда

$$|R(f)| \leq \frac{M_2}{24} h^2 (b - a), \quad (5.6)$$

где  $M_2 = \max_{a \leq x \leq b} |f''(x)|$ .

4. Формула трапеций:

$$\int_a^b f(x) dx \approx \sum_{i=1}^N \frac{f(x_{i-1}) + f(x_i)}{2} h_i. \quad (5.7)$$

Используя (4.3), получим

$$R(f) = -\frac{1}{12} \sum_{i=1}^N f''(\xi_i) h_i^3,$$

$\xi_i \in (x_{i-1}, x_i)$ ,  $i = 1, \dots, N$ , откуда

$$|R(f)| \leq \frac{M_2}{12} h^2 (b - a). \quad (5.8)$$

5. Формула Симпсона:

$$\int_a^b f(x) dx \approx \sum_{i=1}^N \frac{f(x_{i-1}) + 4f(x_{i-1/2}) + f(x_i)}{6} h_i. \quad (5.9)$$

Используя (4.5), получим

$$R(f) = -\frac{1}{2880} \sum_{i=1}^N f^{IV}(\xi_i) h_i^5,$$

$\xi_i \in (x_{i-1}, x_i)$ ,  $i = 1, \dots, N$ , откуда

$$|R(f)| \leq \frac{M_4}{2880} h^4 (b - a), \quad (5.10)$$

где  $M_4 = \max_{a \leq x \leq b} |f^{IV}(x)|$ .

**6. Квадратурные формулы типа Гаусса.** Интерполяционная квадратурная формула с  $n$  узлами при любом расположении узлов точна на любом полиноме степени  $n - 1$ . Однако при удачном расположении узлов точность может повышаться. Так, квадратурная формула центральных прямоугольников оказывается точной на любом полиноме первой степени, а квадратурная формула Симпсона — на любом полиноме третьей степени (см. соответствующие представления остаточных членов (4.4), (4.5)). В связи с этим возникает задача: указать такой способ расположения узлов интерполяционной квадратурной формулы, чтобы сделать ее точной на полиномах возможно более высокой степени. Понятно, что, имея в распоряжении  $n$  свободных параметров, а именно  $x_1, x_2, \dots, x_n$ , можно попытаться построить формулу, точную на любом полиноме степени  $2n - 1$ . Следующая теорема показывает, что большей точности добиться нельзя.

**Теорема 6.1.** *Ни при каком расположении узлов квадратурная формула вида (1.1) не может быть точной на любом полиноме степени  $2n$ .*

ДОКАЗАТЕЛЬСТВО. Пусть  $P_{2n}(x) = (x - x_1)^2(x - x_2)^2 \cdots (x - x_n)^2$ . Тогда

$$\int_a^b p(x)P_{2n}(x) dx > 0, \quad \text{а} \quad \sum_{k=1}^n c_k P_{2n}(x_k) = 0. \quad \square$$

Следующая теорема гарантирует существование квадратурной формулы с  $n$  узлами, точной на любом полиноме степени  $2n - 1$ , и, более того, указывает способ ее построения.

**Теорема 6.2.** *Для того, чтобы квадратурная формула*

$$I(f) = \int_a^b p(x)f(x) dx \approx S(f) = \sum_{k=1}^n c_k f(x_k) \quad (6.1)$$

*была точна на любом полиноме степени  $2n - 1$ , необходимо и достаточно, чтобы:*

- 1) *квадратурная формула была интерполяционной,*
- 2) *узлы квадратурной формулы совпадали с корнями полинома степени  $n$ , ортогонального с весом  $p(x)$  на отрезке  $[a, b]$ .*

ДОКАЗАТЕЛЬСТВО. Пусть выполнены условия 1), 2). Покажем, что квадратурная формула (6.1) точна на любом полиноме  $P_{2n-1}(x)$  степени  $2n - 1$ . Представим полином  $P_{2n-1}(x)$  в виде

$$P_{2n-1}(x) = \omega_n(x)q_{n-1}(x) + r_{n-1}(x), \quad (6.2)$$

где  $\omega_n(x) = (x-x_1)(x-x_2)\cdots(x-x_n)$ , а  $q_{n-1}(x)$ ,  $r_{n-1}(x)$  — полиномы степени  $n-1$ , однозначно определяемые по полиному  $P_{2n-1}(x)$ . По построению полином  $\omega_n(x)$  ортогонален с весом  $p(x)$  на отрезке  $[a, b]$ , поэтому

$$\int_a^b p(x)P_{2n-1}(x) dx = \int_a^b p(x)r_{n-1}(x) dx.$$

Формула (6.1) интерполяционная, поэтому она точна на любом полиноме степени  $n-1$ , следовательно,

$$\int_a^b p(x)r_{n-1}(x) dx = \sum_{k=1}^n c_k r_{n-1}(x_k),$$

но, учитывая (6.2), имеем  $r_{n-1}(x_k) = P_{2n-1}(x_k)$ ,  $k = 1, 2, \dots, n$ . Таким образом,

$$\int_a^b p(x)P_{2n-1}(x) dx = \sum_{k=1}^n c_k P_{2n-1}(x_k).$$

Обратно, если квадратурная формула точна на любом полиноме степени  $2n-1$ , то она точна и на любом полиноме степени  $n-1$ . Осталось показать, что узлы такой квадратурной формулы совпадают с корнями полинома степени  $n$ , ортогонального с весом  $p(x)$  на отрезке  $[a, b]$ . Действительно, полином  $\omega_n(x)q_{n-1}(x)$ , где  $q_{n-1}(x)$  — произвольный полином степени  $n-1$ , есть полином степени  $2n-1$ , следовательно,

$$\int_a^b p(x)\omega_n(x)q_{n-1}(x) dx = \sum_{k=1}^n c_k \omega_n(x_k)q_{n-1}(x_k) = 0,$$

так как  $\omega_n(x_k) = 0$ ,  $k = 1, 2, \dots, n$ .  $\square$

Таким образом, построение квадратурной формулы, точной на любом полиноме степени  $2n-1$ , сводится к следующему.

Сначала вычисляются корни  $x_1, x_2, \dots, x_n$  полинома степени  $n$ , ортогонального с весом  $p(x)$  на отрезке  $[a, b]$ . Затем коэффициенты  $c_k$  вычисляются по формулам

$$c_k = \int_a^b p(x)\varphi_k(x) dx, \quad k = 1, 2, \dots, n,$$



где  $\varphi_k(x)$  — базисные функции Лагранжа, построенные по узлам  $x_1, x_2, \dots, x_n$ .

Квадратурные формулы, точные на любом полиноме степени  $2n-1$ , называют квадратурными формулами наивысшей алгебраической степени точности, или квадратурными формулами типа Гаусса.

Заметим, что из доказательства теоремы 6.2 и свойства единственности ортогональных полиномов вытекает единственность квадратурной формулы наивысшей алгебраической степени точности.

**Теорема 6.3.** *Коэффициенты квадратурной формулы типа Гаусса положительны.*

ДОКАЗАТЕЛЬСТВО. Пусть  $1 \leq k \leq n$ . Положим

$$\psi_k(x) = (x - x_1)^2(x - x_2)^2 \cdots (x - x_{k-1})^2(x - x_{k+1})^2 \cdots (x - x_n)^2.$$

Ясно, что  $\psi_k(x) \geq 0$  — полином степени  $2n-2$ , причем  $\psi_k(x_j) = 0$  при  $j \neq k$ , следовательно,

$$\int_a^b p(x)\psi_k(x) dx = \sum_{j=1}^n c_j\psi_k(x_j) = c_k\psi_k(x_k),$$

т. е.  $c_k > 0$ , так как  $\psi_k(x_k) > 0$ .  $\square$

Из этой теоремы и теоремы 2.1 вытекает, что квадратурные формулы типа Гаусса устойчивы.

Получим представление для остаточного члена квадратурной формулы типа Гаусса.

**Теорема 6.4.** *Пусть функция  $f$  непрерывно дифференцируема  $2n$  раз на отрезке  $[a, b]$ . Тогда для остаточного члена квадратурной формулы типа Гаусса справедливо представление*

$$R(f) = \frac{f^{(2n)}(\xi)}{(2n)!} \int_a^b p(x)\omega_n^2(x) dx, \quad \xi \in (a, b). \quad (6.3)$$

ДОКАЗАТЕЛЬСТВО. Построим полином  $L_{2n-1}(x)$  степени  $2n-1$ , удовлетворяющий условиям:

$$\begin{aligned} L_{2n-1}(x_1) &= f(x_1), \quad L_{2n-1}(x_2) = f(x_2), \quad \dots, \quad L_{2n-1}(x_n) = f(x_n), \\ L'_{2n-1}(x_1) &= f'(x_1), \quad L'_{2n-1}(x_2) = f'(x_2), \quad \dots, \quad L'_{2n-1}(x_n) = f'(x_n), \end{aligned}$$

а также полином  $\tilde{L}_{n-1}$  степени  $n-1$ , удовлетворяющий условиям:

$$\tilde{L}_{n-1}(x_1) = f(x_1), \quad \tilde{L}_{n-1}(x_2) = f(x_2), \quad \dots, \quad \tilde{L}_{n-1}(x_n) = f(x_n).$$

По определению

$$R(f) = \int_a^b p(x)(f(x) - \tilde{L}_{n-1}(x)) dx.$$

Выполним очевидное преобразование

$$R(f) = \int_a^b p(x)(f(x) - L_{2n-1}(x) + L_{2n-1}(x) - \tilde{L}_{n-1}(x)) dx.$$

Вновь используя определение квадратурной формулы типа Гаусса, получим

$$\int_a^b p(x)\tilde{L}_{n-1}(x) = \sum_{k=1}^n c_k f(x_k),$$

а поскольку эта формула точна на любом полиноме степени  $2n - 1$ , то

$$\int_a^b p(x)L_{2n-1}(x) dx = \sum_{k=1}^n c_k L_{2n-1}(x_k),$$

но  $L_{2n-1}(x_k) = f(x_k)$ ,  $k = 1, 2, \dots, n$ , поэтому

$$R(f) = \int_a^b p(x)(f(x) - L_{2n-1}(x)) dx.$$

Используя теперь формулу (6.2), § 1, получим:

$$f(x) - L_{2n-1}(x) = \frac{f^{(2n)}(\zeta(x))}{(2n)!} \omega_n^2(x).$$

Применяя правило Лопиталья, нетрудно убедиться, что функция  $f^{(2n)}(\zeta(x))$  непрерывна по  $x$  на отрезке  $[a, b]$ . Для завершения доказательства теоремы достаточно учесть неотрицательность функции  $p(x)\omega_n^2(x)$  на отрезке  $[a, b]$  и применить теорему о среднем.  $\square$

Приведем примеры квадратурных формул типа Гаусса.

1. Формула Эрмита. Так называют формулу типа Гаусса при  $a = -1$ ,  $b = 1$ ,  $p(x) = 1/\sqrt{1-x^2}$ . Полиномы, ортогональные с указанным весом, — полиномы Чебышева  $T_n(x)$ . Поэтому узлы квадратурной формулы Эрмита вычисляются по явным формулам:

$$x_k = \cos \frac{(2k-1)\pi}{2n}, \quad k = 1, 2, \dots, n.$$

Коэффициенты квадратурной формулы Эрмита тоже удается выразить в явном виде:  $c_k = \pi/n$ ,  $k = 1, 2, \dots, n$ . Отметим, что квадратурная формула Эрмита — это единственная квадратурная формула наивысшей алгебраической степени точности, у которой при любом  $n$  все коэффициенты равны между собой.

2. Квадратурная формула Гаусса (часто называемая также квадратурной формулой Гаусса — Лежандра). В этом случае  $a = -1$ ,  $b = 1$ ,  $p(x) \equiv 1$ . Соответствующие ортогональные полиномы — полиномы Лежандра  $L_n(x)$ . Их нули в явном виде удается найти лишь при малых значениях  $n$ .

Пользуясь формулой Родрига (2.5), § 1, легко подсчитать, что (с точностью до постоянного множителя):

$$L_1(x) = x, \quad L_2(x) = 3x^2 - 1, \quad L_3(x) = 5x^3 - 3x.$$

Теперь нетрудно вычислить узлы и коэффициенты квадратурной формулы Гаусса при соответствующих значениях  $n$ :

- 1)  $n = 1$ ,  $x_1 = 0$ ,  $c_1 = 2$ ;
- 2)  $n = 2$ ,  $x_1 = -1/\sqrt{3}$ ,  $x_2 = 1/\sqrt{3}$ ,  $c_1 = c_2 = 1$ ;
- 3)  $n = 3$ ,  $x_1 = -\sqrt{3/5}$ ,  $x_2 = 0$ ,  $x_3 = \sqrt{3/5}$ ,  $c_1 = c_3 = 5/9$ ,  $c_2 = 8/9$ .

При больших  $n$  узлы и коэффициенты формулы Гаусса вычисляются приближенно. Известны таблицы (см., например, [3]), содержащие соответствующие значения.

---

---

## ГЛАВА 3

# Численные методы решения дифференциальных уравнений

### § 1. Численные методы решения задачи Коши

В этом параграфе будут рассмотрены методы приближенного решения задачи Коши для обыкновенного дифференциального уравнения первого порядка

$$u'(x) = f(x, u), \quad x > x_0, \quad (1)$$

$$u(x_0) = u_0, \quad (2)$$

При построении приближенных методов вводится сетка, т. е. дискретное множество точек  $x_0 < x_1 < x_2 < \dots$  на полуоси  $x \geq x_0$ . В качестве приближенного решения разыскивается сеточная функция  $y(x)$ , т. е. функция, определяемая совокупностью значений  $y(x_0), y(x_1), \dots$ . Для простоты изложения в дальнейшем будем ограничиваться лишь случаем равномерной сетки:  $x_i = x_0 + ih$ ,  $h$  — шаг сетки. Будем также использовать обозначение:  $y_k = y(x_k)$ ,  $k = 0, 1, \dots$ .

**1. Метод, основанный на формуле Тейлора.** Предполагая, что решение задачи (1), (2), т. е. функция  $u(x)$ , имеет  $n + 1$  непрерывную производную, можем написать, что

$$\begin{aligned} u(x_1) = u(x_0 + h) = u(x_0) + hu'(x_0) + \frac{h^2}{2!}u''(x_0) + \frac{h^3}{3!}u'''(x_0) + \dots + \\ + \frac{h^n}{n!}u^{(n)}(x_0) + \frac{h^{n+1}}{(n+1)!}u^{(n+1)}(\xi), \quad \xi \in (x_0, x_1). \end{aligned} \quad (1.1)$$

Покажем, что если функция двух переменных  $f(x, p)$  дифференцируема достаточное число раз, то значения производных функции  $u$  в точке  $x_0$  могут быть вычислены по исходным данным задачи. Действительно, вследствие уравнения (1) и начального условия (2) имеем

$$u'(x_0) = f(x_0, u(x_0)) = f(x_0, u_0) \equiv L_1(x_0, u_0).$$

Далее,

$$u''(x) = (f(x, u))' = f_x(x, u) + f_u(x, u)u' = f_x(x, u) + f_u(x, u)f(x, u),$$

следовательно,

$$u''(x_0) = f_x(x_0, u_0) + f_u(x_0, u_0)L_1(x_0, u_0) \equiv L_2(x_0, u_0).$$

Точно так же

$$\begin{aligned} u'''(x) &= (f_x(x, u) + f_u(x, u)f(x, u))' = \\ &= f_{xx}(x, u) + 2f_{xu}(x, u)f(x, u) + f_{uu}(x, u)(f(x, u))^2 + f_u(x, u)u''(x), \end{aligned}$$

поэтому

$$\begin{aligned} u'''(x_0) &= f_{xx}(x_0, u_0) + 2f_{xu}(x_0, u_0)L_1(x_0, u_0) + \\ &+ f_{uu}(x_0, u_0)(L_1(x_0, u_0))^2 + f_u(x_0, u_0)L_2(x_0, u_0) \equiv L_3(x_0, u_0). \end{aligned}$$

Аналогичные выражения можно получить для производной любого порядка:  $u^{(k)}(x_0) = L_k(x_0, u_0)$ .

Таким образом,

$$\begin{aligned} u(x_1) &= u(x_0) + hL_1(x_0, u_0) + \frac{h^2}{2!}L_2(x_0, u_0) + \dots + \frac{h^n}{n!}L_n(x_0, u_0) + \\ &+ \frac{h^{n+1}}{(n+1)!}u^{(n+1)}(\xi), \quad \xi \in (x_0, x_1). \end{aligned}$$

Отбрасывая в этом равенстве остаточный член формулы Тейлора и обозначая получаемое приближение к  $u(x_1)$  через  $y(x_1)$ , можем написать

$$y(x_1) = u(x_0) + hL_1(x_0, u_0) + \frac{h^2}{2!}L_2(x_0, u_0) + \dots + \frac{h^n}{n!}L_n(x_0, u_0).$$

Теперь естественно считать, что вообще

$$y_{i+1} = y_i + hF(x_i, y_i), \quad i = 0, 1, \dots, \quad y_0 = u_0. \quad (1.2)$$

Здесь

$$F(x_i, y_i) = L_1(x_i, y_i) + \frac{h}{2!}L_2(x_i, y_i) + \dots + \frac{h^{(n-1)}}{n!}L_n(x_i, y_i). \quad (1.3)$$

Таким образом, построен метод вычисления приближенного решения задачи (1), (2). При реализации метода (1.2), (1.3) необходимо вычислять частные производные функции  $f(x, p)$ , причем с увеличением точности метода, т. е. с увеличением  $n$ , порядок требуемых производных возрастает, а формулы, по которым вычисляются выражения  $L_j(x_i, y_i)$ , становятся все более громоздкими. Поэтому метод,

основанный на формуле Тейлора, довольно редко используется на практике. Правда, в последнее время в связи с появлением эффективных компьютерных программ аналитических вычислений (например, Maple) отношение к этому методу меняется. Тем не менее, наиболее распространены такие приближенные методы решения задачи Коши, которые требуют для своей реализации умения вычислять лишь значения функции  $f(x, p)$  по заданным значениям аргументов.

**2. Методы типа Рунге — Кутта.** Пусть  $u(x)$  — решение задачи (1), (2). Проинтегрировав равенство (1) по отрезку  $[x_0, x_1]$ , получим

$$u(x_1) = u(x_0) + \int_{x_0}^{x_1} f(x, u(x)) dx. \quad (2.1)$$

Дальнейшее, фактически, основано на том или ином способе приближенного вычисления участвующего здесь интеграла.

1. Используем сначала самую простую формулу — формулу левых прямоугольников. Получим

$$u(x_1) \approx u_0 + hf(x_0, u_0).$$

Это приводит к следующей расчетной формуле:

$$y_1 = y_0 + hf(x_0, y_0).$$

Вообще,

$$y_{i+1} = y_i + hf(x_i, y_i), \quad i = 0, 1, 2, \dots, \quad y_0 = u_0. \quad (2.2)$$

Метод (2.2) называется явным методом Эйлера (чаще, просто методом Эйлера) или методом ломаных<sup>1)</sup>.

Это — самый простой вариант метода Рунге — Кутта. Метод имеет простую и полезную геометрическую интерпретацию. Из точки  $(x_0, u_0)$  выпускается прямая по касательной к интегральной кривой уравнения (1), проходящей через точку  $(x_0, u_0)$ . Из точки  $(x_1, y_1)$  прямая идет уже по касательной к интегральной кривой, проходящей через точку  $(x_1, y_1)$ , и так далее. В результате получается ломаная, аппроксимирующая искомую интегральную кривую. Понятно, что этот метод не может обеспечить высокой точности, если шаг сетки не слишком мал, тем не менее он применяется довольно часто.

2. Неявный метод Эйлера. Используем для приближенного вычисления интеграла формулу правых прямоугольников

$$u(x_1) \approx u_0 + hf(x_1, u_1).$$

<sup>1)</sup>Он совпадает с методом, основанным на формуле Тейлора, при  $n = 1$ .

В результате приходим к следующему методу:

$$y_{i+1} = y_i + hf(x_{i+1}, y_{i+1}), \quad i = 0, 1, \dots, \quad y_0 = u_0.$$

Этот метод называется неявным методом Эйлера. Значение  $y_{i+1}$  не может быть явно определено по известному  $y_i$ , а разыскивается как решение уравнения

$$y_{i+1} - hf(x_{i+1}, y_{i+1}) = y_i. \quad (2.3)$$

При этом применяются приближенные (итерационные) методы. Как правило, используется метод простой итерации, состоящий в следующем. Выбирается некоторое начальное приближение  $y_{i+1}^{(0)}$  к  $y_{i+1}$ . Обычно полагают  $y_{i+1}^{(0)} = y_i$  или определяют  $y_{i+1}^{(0)}$  при помощи явного метода Эйлера по известному  $y_i$ . Затем строят последовательность приближений

$$y_{i+1}^{(k+1)} = y_i + hf(x_{i+1}, y_{i+1}^{(k)}), \quad k = 0, 1, \dots$$

Этот итерационный метод быстро сходится, если  $h$  достаточно мало. На практике редко выполняют больше двух-трех итераций.

3. Формула центральных прямоугольников (метод предиктор-корректор). Увеличим точность приближения интеграла. Для этого используем формулу центральных прямоугольников:

$$u(x_1) \approx u(x_0) + hf(x_{1/2}, u(x_{1/2})),$$

где  $x_{1/2} = (x_0 + x_1)/2$ . Значение  $u(x_{1/2})$  приближенно определяется при помощи явного метода Эйлера

$$u(x_{1/2}) \approx u(x_0) + \frac{h}{2}f(x_0, u_0).$$

Расчетные формулы таковы:

$$y_{i+1} = y_i + hf(x_{i+1/2}, y_{i+1/2}), \quad i = 1, 2, \dots, \quad (2.4)$$

где

$$y_{i+1/2} = y_i + \frac{h}{2}f(x_i, y_i), \quad x_{i+1/2} = (x_i + x_{i+1})/2. \quad (2.5)$$

Все вычисления проводятся по явным формулам. Шаг (2.5) называется предиктором, шаг (2.4) — корректором, то есть сначала по формуле (2.5) мы выполняем как бы предсказание (to predict — предсказывать), а затем по формуле (2.4) уточнение (to correct — исправлять, уточнять) значения  $y_{i+1}$ . Здесь в отличие от явного метода переход от  $y_i$  к  $y_{i+1}$  требует вычисления двух значений функции  $f$ .

4. Правило трапеций. Используем для приближенного вычисления интеграла формулу трапеций:

$$u(x_1) \approx u(x_0) + \frac{h}{2} (f(x_0, u_0) + f(x_1, u(x_1))).$$

Отсюда вытекает расчетная формула:

$$y_{i+1} = y_i + \frac{h}{2} (f(x_i, y_i) + f(x_{i+1}, y_{i+1})), \quad i = 0, 1, 2, \dots, \quad y_0 = u_0. \quad (2.6)$$

Это — неявный метод. По поводу его реализации можно сказать то же, что и по поводу реализации неявного метода Эйлера.

Комбинируя правило трапеций с явной формулой Эйлера, получим еще один вариант метода предиктор-корректор:

$$y_{i+1} = y_i + \frac{h}{2} (f(x_i, y_i) + f(x_{i+1}, \bar{y}_{i+1})), \quad i = 0, 1, 2, \dots, \quad y_0 = u_0, \quad (2.7)$$

где

$$\bar{y}_{i+1} = y_i + hf(x_i, y_i). \quad (2.8)$$

Сначала выполняется предсказание по формуле (2.8), а затем коррекция по формуле (2.7).

5. Метод Рунге — Кутта четвертого порядка точности. В основе этого метода лежит квадратурная формула Симпсона. Используются также предсказания по формуле Эйлера. Не останавливаясь на подробном выводе, дадим окончательные расчетные формулы:

$$y_{i+1} = y_i + \frac{1}{6} (k_1 + 2k_2 + 2k_3 + k_4), \quad i = 0, 1, 2, \dots, \quad y_0 = u_0, \quad (2.9)$$

где

$$\begin{aligned} k_1 &= hf(x_i, y_i), & k_2 &= hf\left(x_i + \frac{h}{2}, y_i + \frac{k_1}{2}\right), \\ k_3 &= hf\left(x_i + \frac{h}{2}, y_i + \frac{k_2}{2}\right), & k_4 &= hf(x_i + h, y_i + k_3). \end{aligned}$$

Этот метод (с небольшими модификациями) на практике применяется, по-видимому, наиболее часто.

**3. Элементы теории одношаговых методов решения задачи Коши.** Метод, основанный на формуле Тейлора, а также методы типа Рунге — Кутта могут быть представлены в общем виде

$$\frac{y_{i+1} - y_i}{h} = F(x_i, h, y_i, y_{i+1}), \quad i = 0, 1, \dots, \quad y_0 = u_0. \quad (3.1)$$



Здесь  $F(x_i, h, y_i, y_{i+1})$  — некоторая функция, выбор которой и определяет тот или иной метод. Например, для неявного метода типа Рунге — Кутты, основанного на формуле трапеций,

$$F(x_i, h, y_i, y_{i+1}) = \frac{1}{2} (f(x_i, y_i) + f(x_i + h, y_{i+1})).$$

Понятно, что значения  $u_i, u_{i+1}$  точного решения задачи Коши не удовлетворяют равенству (3.1). После их подстановки возникает невязка

$$\psi_i = \frac{u_{i+1} - u_i}{h} - F(x_i, h, u_i, u_{i+1}), \quad (3.2)$$

которая называется погрешностью аппроксимации.

Интуитивно ясно, что чем меньше погрешность аппроксимации, тем точнее метод. Ниже будет показано, что при определенных условиях это действительно так.

Исследуем поведение погрешности аппроксимации в зависимости от величины шага сетки для построенных выше методов.

1) Метод, основанный на формуле Тейлора. В этом случае и из способа построения метода сразу вытекает, что  $\psi_i = O(h^n)$ .

2) Явный метод Эйлера. Этот метод есть частный случай метода, основанного на формуле Тейлора, при  $n = 1$ .

3) Неявный метод Эйлера. В этом случае

$$\psi_i = \frac{u_{i+1} - u_i}{h} - f(x_{i+1}, u_{i+1}).$$

Используя оценку погрешности квадратурной формулы правых прямоугольников, можем написать:

$$\begin{aligned} hf(x_{i+1}, u_{i+1}) &= \int_{x_i}^{x_{i+1}} f(x, u(x)) dx + O(h^2) = \int_{x_i}^{x_{i+1}} u'(x) dx + O(h^2) = \\ &= (u(x_{i+1}) - u(x_i)) + O(h^2), \end{aligned}$$

следовательно,  $\psi_i = O(h)$ .

4) Метод предиктор-корректор, основанный на формуле трапеций. В соответствии с формулами (2.4), (2.5)

$$\begin{aligned} \psi_i &= \frac{u_{i+1} - u_i}{h} - f\left(x_{i+1/2}, u(x_i) + \frac{h}{2}f(x_i, u(x_i))\right) = \\ &= \frac{u_{i+1} - u_i}{h} - f\left(x_{i+1/2}, u(x_i) + \frac{h}{2}u'(x_i)\right). \end{aligned}$$



Приведем общий результат о сходимости одношаговых методов вида (3.3). Предположим, что задача (1), (2) решается на отрезке  $x_0 \leq x \leq x_0 + l$ . Пусть  $h = l/N$ ,  $N$  — целое положительное число,  $x_i = x_0 + ih$ ,  $i = 0, \dots, N$ . Будем обозначать через  $y$  приближенное решение, построенное по методу (3.3), через  $u$  — точное решение задачи (1), (2),  $z_i = z(x_i) = y_i - u(x_i)$  — погрешность метода.

**Теорема 3.1.** Пусть функция  $f(x, p)$  непрерывно дифференцируема по  $p$ ,

$$|f_p(x, p)| \leq M = \text{const} \quad \forall x \in [x_0, x_0 + l], \quad p \in (-\infty, \infty), \quad (3.4)$$

погрешность аппроксимации метода есть величина порядка  $h^s$ .

Тогда

$$|z_i| \leq Ch^s, \quad C = \text{const}, \quad i = 0, 1, \dots, N. \quad (3.5)$$

**ДОКАЗАТЕЛЬСТВО.** Покажем, прежде всего, что выполнение условия (3.4) обеспечивает существование такой постоянной  $M_1 > 0$ , что для функции  $F(x, h, y)$  справедлива оценка

$$|F_y(x, h, y)| \leq M_1. \quad (3.6)$$

Для этого установим, что  $|\partial k_r(x, h, y)/\partial y| \leq c_r$  где  $c_r = \text{const}$ ,  $r = 1, 2, \dots, q$ . При  $r = 1$  это непосредственно вытекает из условия (3.4). При  $r = 2$  имеем

$$\begin{aligned} \frac{\partial k_2(x, h, y)}{\partial y} &= \frac{\partial f(x + \alpha_2 h, y + \beta_{21} h k_1(x, h, y))}{\partial y} = \\ &= f_p(x + \alpha_2 h, y + \beta_{21} h k_1(x, h, y))(1 + \beta_{21} h k_{1y}(x, h, y)), \end{aligned}$$

следовательно,  $|\partial k_2(x, h, y)/\partial y| \leq M(1 + \beta_{21} h c_1) \leq M(1 + \beta_{21} l c_1) = c_2$ . Для  $r > 2$  оценки проводятся аналогично. Используя теперь определение погрешности аппроксимации, для любого  $i = 0, 1, \dots, N - 1$  можем написать

$$u_{i+1} = u_i + F(x_i, h, u_i) + h\psi_i. \quad (3.7)$$

Вычитая почленно из равенства (3.3) равенство (3.7), получим

$$z_{i+1} = z_i + h(F(x_i, h, y_i) - F(x_i, h, u_i)) - h\psi_i,$$

откуда вследствие формулы конечных приращений Лагранжа и оценки (3.6) имеем  $|z_{i+1}| \leq (1 + M_1 h)|z_i| + h|\psi_i|$ . Точно так же получаем, что  $|z_i| \leq (1 + M_1 h)|z_{i-1}| + h|\psi_{i-1}|$ , следовательно,

$$|z_{i+1}| \leq (1 + M_1 h)^2 |z_{i-1}| + h(1 + M_1 h)|\psi_{i-1}| + h|\psi_i|.$$

Продолжая аналогичным образом, получим:

$$|z_{i+1}| \leq (1 + M_1 h)^{i+1} |z_0| + h \sum_{k=0}^i (1 + M_1 h)^{i-k} |\psi_k|.$$

Заметим теперь, что для любого целого  $r$ ,  $0 \leq r \leq N$ ,

$$(1 + M_1 h)^r = (1 + M_1 h)^{\frac{x_r - x_0}{h}} = (1 + M_1 h)^{\frac{1}{M_1 h} M_1 (x_r - x_0)} \leq (1 + M_1 h)^{\frac{1}{M_1 h} M_1 l}.$$

Но, как хорошо известно,

$$(1 + v)^{\frac{1}{v}} \leq e \quad \text{при любом } v > 0,$$

поэтому

$$(1 + M_1 h)^{\frac{1}{M_1 h} M_1 l} \leq e^{M_1 l},$$

и, следовательно,

$$|z_{i+1}| \leq e^{M_1 l} |z_0| + e^{M_1 l} \max_{0 \leq k \leq N} |\psi_k| h (i + 1) \leq e^{M_1 l} (|z_0| + l \max_{0 \leq k \leq N} |\psi_k|).$$

Вспоминая теперь, что  $z_0 = y_0 - u_0 = 0$ ,  $|\psi_k| \leq ch^s$ ,  $c = \text{const}$ , получим, что для любого  $i = 0, 1, \dots, N$  справедлива оценка

$$|z_{i+1}| \leq e^{M_1 l} l c h^s. \quad \square$$

**4. Методы типа Адамса.** В отличие от одношаговых методов в методах типа Адамса при вычислении каждого нового значения  $y_{i+1}$  используются не только  $y_i$ , но и значения приближенного решения в нескольких предыдущих точках сетки. Благодаря этому методы типа Адамса оказываются экономичнее одношаговых методов.

Опишем способ построения методов типа Адамса, основанный на использовании интерполяционных многочленов. Предположим, что нам известны значения решения задачи (1), (2) в точках сетки  $x_0, x_1, \dots, x_{k-1}$ . Интегрируя уравнение (1) по отрезку  $[x_{k-1}, x_k]$ , получим

$$u(x_k) = u(x_{k-1}) + \int_{x_{k-1}}^{x_k} f(x, u(x)) dx. \quad (4.1)$$

Проинтерполируем функцию  $f(x, u(x))$  по ее значениям в точках  $x_0, x_1, \dots, x_{k-1}$ , т. е. представим эту функцию в виде

$$f(x, u(x)) = \sum_{l=0}^{k-1} f(x_l, u(x_l)) \varphi_l(x) + R_k(x), \quad (4.2)$$

где  $\varphi_l(x)$ ,  $l = 0, 1, \dots, k-1$  — базисные функции Лагранжа,  $R_k(x)$  — остаточный член интерполяционного полинома. Подставляя выражение (4.2) в равенство (4.1), получим

$$u(x_k) = u(x_{k-1}) + \sum_{l=0}^{k-1} c_l f(x_l, u(x_l)) + \tilde{R}_k, \quad (4.3)$$

где

$$c_l = \int_{x_{k-1}}^{x_k} \varphi_l(x) dx = h \hat{c}_l, \quad \tilde{R}_k = \int_{x_{k-1}}^{x_k} R_k(x) dx. \quad (4.4)$$

Здесь коэффициенты  $\hat{c}_l$ ,  $l = 0, 1, \dots, k-1$ , зависят только от  $k$  и  $l$  (ср. с соответствующими выкладками, выполненными при построении квадратурных формул Ньютона — Котеса). Отбрасывая в равенстве (4.3) остаточный член, естественно считать, что приближенные значения решения задачи в точках  $x_0, x_1, \dots, x_k$ , т. е.  $y_0, y_1, \dots, y_k$  связаны соотношением

$$y_k = y_{k-1} + h \sum_{l=0}^{k-1} \hat{c}_l f(x_l, y_l)$$

и, вообще,

$$y_i = y_{i-1} + h \sum_{l=0}^{k-1} \hat{c}_l f(x_{i-k+l}, y_{i-k+l}), \quad i = k, k+1, \dots \quad (4.5)$$

Если считать значения  $y_0, y_1, \dots, y_{k-1}$  известными, то формулу (4.5) можно рассматривать как рекуррентную для определения всех последующих значений приближенного решения. Значения  $y_1, y_2, \dots, y_{k-1}$  вычисляют, обычно, при помощи метода Рунге — Кутты.

Метод (4.5) называют явным или экстраполяционным методом Адамса. Последнее название объясняется тем, что интерполяционный полином, построенный для функции  $f(x, u(x))$  по узлам  $x_0, x_1, \dots, x_{k-1}$ , экстраполируется на отрезок  $[x_{k-1}, x_k]$ .

Приближая функцию  $f(x, u(x))$  на отрезке  $[x_{k-1}, x_k]$ , можно использовать интерполяционный полином, построенный по узлам  $x_0, x_1, \dots, x_k$ . В результате, вместо (4.5) получаем соотношения

$$y_i = y_{i-1} + h \sum_{l=0}^k \tilde{c}_l f(x_{i-k+l}, y_{i-k+l}), \quad i = k, k+1, \dots, \quad (4.6)$$

определяющее интерполяционный метод Адамса.

Для отыскания  $y_i$  по известным  $y_{i-1}, y_{i-2}, \dots, y_{i-k}$  приходится решать нелинейное, вообще говоря, уравнение

$$y_i - h\tilde{c}_k f(x_i, y_i) = y_{i-1} + h \sum_{l=0}^k \tilde{c}_l f(x_{i-k+l}, y_{i-k+l}). \quad (4.7)$$

Поэтому метод (4.6) часто называют неявным методом Адамса.

Проведем сравнение по трудоемкости методов типа Рунге — Кутта с явным методом Адамса. Ясно, что при вычислении значения  $y_i$  по методу (4.5) приходится вычислять лишь одно значение функции  $f$ , а именно  $f(x_{i-1}, y_{i-1})$ . Все остальные значения функции  $f$ , участвующие в формуле (4.5), уже подсчитаны на предыдущих шагах рекуррентного процесса и могут быть сохранены в памяти компьютера. При применении метода Рунге — Кутта для перехода от  $y_{i-1}$  к  $y_i$  требуется (в зависимости от порядка метода) вычислить несколько значений функции  $f$ . Таким образом, если функция  $f$  достаточно сложная, то методы типа Рунге — Кутта существенно уступают в смысле трудоемкости явному методу Адамса.

Трудоемкость неявного метода Адамса определяется количеством итераций, затрачиваемых для решения нелинейного уравнения (4.7). Обычно с этой целью применяют метод простой итерации. При достаточно малом шаге сетки  $h$  для достижения приемлемой точности чаще всего оказывается достаточным провести две-три итерации и соответственно два-три раза обратиться к процедуре вычисления функции  $f$ . Таким образом, и неявный метод Адамса, как правило, экономичнее методов типа Рунге — Кутта.

Понятие погрешности аппроксимации для методов типа Адамса вводится по аналогии с одношаговыми методами. Например, для явного метода Адамса (4.5) погрешность аппроксимации есть

$$\psi_i = \frac{u_{i+1} - u_i}{h} - \sum_{l=0}^{k-1} \hat{c}_l f(x_{i-k+l}, u_{i-k+l}).$$

Оценки погрешностей аппроксимации методов (4.5), (4.6) нетрудно получить, используя представление остаточного члена интерполяционного полинома.

Можно показать, что если начальные условия определены с точностью порядка погрешности аппроксимации, то есть  $y_i = u_i + O(h^s)$  для  $i = 0, 1, \dots, k-1$ ,  $\psi_i = O(h^s)$ ,  $i = k, k+1, \dots$ , функция  $f(x, p)$  удовлетворяет условию (3.4), то  $|z_i| = O(h^s)$ ,  $i = 0, 1, \dots$ , иными словами, метод имеет точность порядка  $h^s$ .

В заключение приведем примеры методов типа Адамса с указанием порядка погрешности аппроксимации.

1) Явные методы:

$$\frac{y_{i+1} - y_i}{h} = \frac{3}{2}f(x_i, y_i) - \frac{1}{2}f(x_{i-1}, y_{i-1}), \quad s = 2,$$

$$\frac{y_{i+1} - y_i}{h} = \frac{1}{12}(23f(x_i, y_i) - 16f(x_{i-1}, y_{i-1}) + 5f(x_{i-2}, y_{i-2})), \quad s = 3,$$

$$\begin{aligned} \frac{y_{i+1} - y_i}{h} = \frac{1}{24} & (55f(x_i, y_i) - 59f(x_{i-1}, y_{i-1}) + \\ & + 37f(x_{i-2}, y_{i-2}) - 9f(x_{i-3}, y_{i-3})), \quad s = 4. \end{aligned}$$

2) Неявные методы:

$$\frac{y_{i+1} - y_i}{h} = \frac{1}{2}(f(x_{i+1}, y_{i+1}) + f(x_i, y_i)), \quad s = 2,$$

$$\frac{y_{i+1} - y_i}{h} = \frac{1}{12}(5f(x_{i+1}, y_{i+1}) + 8f(x_i, y_i) - f(x_{i-1}, y_{i-1})), \quad s = 3,$$

$$\begin{aligned} \frac{y_{i+1} - y_i}{h} = \frac{1}{24} & (9f(x_{i+1}, y_{i+1}) + 19f(x_i, y_i) - \\ & - 5f(x_{i-1}, y_{i-1}) + f(x_{i-2}, y_{i-2})), \quad s = 4. \end{aligned}$$

## § 2. Методы решения уравнений с частными производными

**1. Сеточные методы решения краевых задач для обыкновенных дифференциальных уравнений.** Рассматривается обыкновенное дифференциальное уравнение второго порядка

$$-(pu')' + qu = f, \quad 0 < x < l, \quad (1.1)$$

с граничными условиями

$$u(0) = u(1) = 0, \quad (1.2)$$

где,  $p(x)$ ,  $q(x)$ ,  $f(x)$  — заданные функции, причем

$$p(x) \geq c_0 = \text{const} > 0, \quad q(x) \geq 0. \quad (1.3)$$

Задача (1.1), (1.2) может быть, например, интерпретирована как задача о стационарном распределении температуры в стержне, на боковой поверхности которого происходит теплообмен по закону Ньютона со средой, имеющей заданную температуру. Наш интерес к этой задаче обусловлен не столько важностью соответствующих приложений, сколько тем, что задача (1.1), (1.2) может рассматриваться как хорошая модель для отработки методов численного решения более сложных двумерных и трехмерных эллиптических уравнений.

Построим на отрезке  $[0, l]$  равномерную сетку

$$\omega = \{x_i = ih, i = 0, 1, 2, \dots, N, nh = 1\}$$

с шагом  $h$ . Если  $u(x)$  — дифференцируемая функция, то по определению производной

$$u'(x_i) \approx \frac{u(x_i + h) - u(x_i)}{h} = \frac{u_{i+1} - u_i}{h}$$

при достаточно малом  $h$ . Выражение, стоящее в правой части этого приближенного равенства, называется разностным отношением. Понятно, что это не единственный способ приближения производной. Точно так же

$$u'(x_i) \approx \frac{u(x_i) - u(x_i - h)}{h} = \frac{u_i - u_{i-1}}{h}.$$

Приняты следующие наименования и обозначения:

$$\frac{u_{i+1} - u_i}{h} = u_{x,i} \quad \text{— разностное отношение вперед,}$$

или правое разностное отношение,

$$\frac{u_i - u_{i-1}}{h} = u_{\bar{x},i} \quad \text{— разностное отношение назад,}$$

или левое разностное отношение.

Вторую производную  $u''(x_i)$  естественно приближать выражением

$$u_{\bar{x}x,i} = (u_{\bar{x}})_{x,i} = \frac{1}{h} \left( \frac{u_{i+1} - u_i}{h} - \frac{u_i - u_{i-1}}{h} \right) = \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2};$$

$u_{\bar{x}x,i}$  называется второй разделенной разностью или вторым разностным отношением.

Более сложное дифференциальное выражение  $(pu')'$  будем приближать — разностным  $(au_{\bar{x}})_x$ , где  $a(x)$  — сеточная функция, в определенном смысле близкая  $p(x)$ . Как именно следует выбирать  $a(x)$ ,



будет сказано чуть позже. Таким образом, уравнению (1.1) естественно поставить в соответствие приближенное равенство

$$-(au_{\bar{x}})_x + bu \approx \varphi, \quad x \in \overset{\circ}{\omega},$$

где  $\overset{\circ}{\omega} = \{x_i = ih, i = 1, 2, \dots, N - 1\}$  — множество внутренних точек сетки,  $a(x)$ ,  $b(x)$ ,  $\varphi(x)$  — сеточные функции, вычисляемые некоторым образом по известным функциям  $p(x)$ ,  $q(x)$ ,  $f(x)$ . Чаще всего полагают

$$a(x) = p(x - h/2), \quad b(x) = q(x), \quad \varphi(x) = f(x). \quad (1.4)$$

Определим теперь сеточную функцию  $y(x)$  как решение системы линейных алгебраических уравнений:

$$-(ay_{\bar{x}})_x + by = \varphi, \quad x \in \overset{\circ}{\omega}. \quad (1.5)$$

Приведем более подробную запись этой системы:

$$-\frac{a_{i+1}(y_{i+1} - y_i) - a_i(y_i - y_{i-1}))}{h^2} + b_i y_i = \varphi_i, \quad i = 1, 2, \dots, N - 1,$$

или

$$A_i y_{i+1} - B_i y_i + C_i y_{i-1} = -\varphi_i, \quad i = 1, 2, \dots, N - 1, \quad (1.6)$$

где

$$A_i = a_{i+1}/h^2, \quad C_i = a_i/h^2, \quad B_i = b_i + (a_{i+1} + a_i)/h^2. \quad (1.7)$$

Ясно, что система (1.5) содержит уравнений на два меньше, чем неизвестных. Для того, чтобы пополнить ее, воспользуемся граничными условиями (1.2) и положим

$$y_0 = y_N = 0. \quad (1.8)$$

Совокупность уравнений (1.5), (1.8) образует полную систему линейных алгебраических уравнений относительно  $y_0, y_1, \dots, y_N$  и называется разностной схемой для задачи (1.1), (1.2). Для решения этой трехдиагональной системы уравнений чаще всего применяют метод прогонки (см. с. 16). Условия (1.3) влекут выполнение неравенств

$$|A_i| + |C_i| < |B_i|, \quad i = 1, 2, \dots, N - 1. \quad (1.9)$$

Именно они обеспечивают реализуемость метода прогонки.

Займемся теперь исследованием погрешности аппроксимации и построением функций  $a(x)$ ,  $b(x)$ ,  $\varphi(x)$ . Как и ранее, под погрешностью аппроксимации будем понимать невязку

$$\psi_i = -(au_{\bar{x}})_{x,i} + bu_i - \varphi_i,$$

возникающую при подстановке в уравнение (1.5) точного решения задачи (1.1), (1.2). Более подробно:

$$\psi_i = -\frac{a_{i+1}(u_{i+1} - u_i) - a_i(u_i - u_{i-1})}{h^2} + b_i y_i - \varphi_i.$$

Воспользуемся формулой Тейлора:

$$u_{i+1} = u\left(x_{i+1/2} + \frac{h}{2}\right) = u_{i+1/2} + \frac{h}{2}u'_{i+1/2} + \frac{1}{2!}\left(\frac{h}{2}\right)^2 u''_{i+1/2} + \\ + \frac{1}{3!}\left(\frac{h}{2}\right)^3 u'''_{i+1/2} + \dots,$$

$$u_i = u\left(x_{i+1/2} - \frac{h}{2}\right) = u_{i+1/2} - \frac{h}{2}u'_{i+1/2} + \frac{1}{2!}\left(\frac{h}{2}\right)^2 u''_{i+1/2} - \\ - \frac{1}{3!}\left(\frac{h}{2}\right)^3 u'''_{i+1/2} + \dots,$$

следовательно,

$$\frac{u_{i+1} - u_i}{h} = u'_{i+1/2} + \frac{2}{3!}\frac{h^2}{2^3}u'''_{i+1/2} + O(h^4).$$

Аналогично,

$$\frac{u_i - u_{i-1}}{h} = u'_{i-1/2} + \frac{2}{3!}\frac{h^2}{2^3}u'''_{i-1/2} + O(h^4).$$

Отсюда с очевидностью вытекает, что если выбрать  $a$ ,  $b$ ,  $\varphi$  в соответствии с (1.4), то, с одной стороны, условия (1.9) будут выполнены, а, с другой — погрешность аппроксимации будет иметь второй порядок малости относительно шага сетки:

$$\psi_i = O(h^2). \quad (1.10)$$

Ясно, конечно, что при этом необходимо, чтобы функции  $u(x)$ ,  $p(x)$  имели достаточное число непрерывных производных.

Нужно отметить, что коэффициенты разностной схемы можно вычислять и другими способами. Например, часто используют соотношения

$$a_i = \frac{p_i + p_{i-1}}{2}, \text{ или } a_i = \left[ \frac{1}{h} \int_{x_{i-1}}^{x_i} \frac{1}{p(\xi)} d\xi \right]^{-1}.$$

При этом оценка (1.10) сохраняется.

**2. Вариационные методы.** Наряду с разностными методами для решения граничных задач для эллиптических уравнений широко применяются вариационные методы. Между разностными и вариационными методами существует довольно тесная связь, наиболее ярко проявляющаяся в методе конечных элементов. Вновь мы ограничимся простейшей эллиптической граничной задачей (1.1), (1.2). Рассмотрим наряду с этой задачей так называемый энергетический функционал (функционал Лагранжа)

$$F(u) = \frac{1}{2} \int_0^l (pu'^2 + qu^2) dx - \int_0^l fudx. \quad (2.1)$$

Отметим, что если интерпретировать задачу (1.1), (1.2) как задачу о равновесии струны ( $q(x)$  можно при этом трактовать как коэффициент жесткости упругого основания (постели), на котором находится струна), то  $F(u)$  — потенциальная энергия системы струна — внешние силы.

**Теорема 2.1.** Пусть функция  $u(x)$  доставляет минимальное значение функционалу  $F$  на множестве функций, удовлетворяющих граничным условиям (1.2). Тогда  $u(x)$  — решение задачи (1.1), (1.2).

**ДОКАЗАТЕЛЬСТВО.** Рассмотрим наряду с функцией  $u(x)$  функцию  $u(x) + t\eta(x)$ , где  $\eta(x)$  удовлетворяет граничным условиям (1.2), а  $t$  — вещественное число. Ясно, что поскольку на функции  $u(x)$  функционал достигает минимального значения, то

$$F(u + t\eta) \geq F(u), \quad (2.2)$$

значит, функция вещественного переменного

$$\varphi(t) = F(u + t\eta)$$

достигает минимального значения при  $t = 0$ , следовательно,

$$\varphi'(0) = 0. \quad (2.3)$$

Подсчитаем  $\varphi'(0)$ . Для этого запишем  $F(u + t\eta)$  более подробно. Имеем:

$$F(u + t\eta) = \frac{1}{2} \int_0^l (pu'^2 + qu^2) dx - \int_0^l f u dx + t \int_0^l (pu'\eta' + qu\eta) dx - \\ - t \int_0^l f \eta dx + \frac{t^2}{2} \int_0^l (p\eta'^2 + q\eta^2),$$

следовательно,

$$\varphi'(0) = \int_0^l (pu'\eta' + qu\eta) dx - \int_0^l f \eta dx,$$

и условие (2.3) принимает вид:

$$\int_0^l (pu'\eta' + qu\eta) dx = \int_0^l f \eta dx. \quad (2.4)$$

Уравнение (2.4) часто называют вариационным уравнением, соответствующим функционалу  $F(u)$ . Оно представляет собой необходимое условие минимума функционала.

Покажем, что уравнение (1.1) вытекает из (2.4). Действительно, используя формулу интегрирования по частям и учитывая граничные условия для функции  $\eta(x)$ , получим

$$\int_0^l pu'\eta' dx = - \int_0^l (pu')' \eta dx + pu'\eta \Big|_0^l = - \int_0^l (pu')' dx,$$

то есть

$$\int_0^l \left( -(pu')' + qu - f \right) \eta dx = 0.$$

Используя теперь то, что функция  $\eta$  в равенстве (2.4) произвольна, покажем, что функция  $\zeta(x) \equiv -(pu')' + qu - f \equiv 0$  на интервале  $(0, l)$ . Действительно, если предположить противное, то можно указать такую точку  $x_0 \in (0, l)$ , что  $\zeta(x_0) \neq 0$ . Для определенности можно считать, что  $\zeta(x_0) > 0$ . Случай противоположного знака исследуется

точно так же. Вследствие непрерывности  $\zeta(x)$  найдется окрестность  $(x_0 - \varepsilon, x_0 + \varepsilon) \subset (0, l)$ , на которой функция  $\zeta(x)$  сохраняет знак. Выберем теперь функцию  $\eta(x)$  так, чтобы она была положительной внутри указанной окрестности и равна тождественно нулю в остальных точках интервала  $(0, l)$ . При таком выборе функции  $\eta(x)$  интеграл в левой части равенства есть

$$\int_{x_0 - \varepsilon}^{x_0 + \varepsilon} \zeta(x)\eta(x)dx > 0,$$

что невозможно. Остается принять, что  $\zeta(x) \equiv 0$ , т. е. уравнение (1.1) выполнено.  $\square$

Уравнение (1.1), возникающее как необходимое условие минимума функционала  $F(u)$ , часто называют уравнением Эйлера.

Доказанная нами теорема позволяет заменить задачу (1.1), (1.2) задачей минимизации функционала  $F$  на множестве функций, удовлетворяющих граничным условиям (1.2). Основанные на такой замене приближенные методы решения задачи (1.1), (1.2) называют вариационными.

**3. Метод Ритца.** Это — исторически первый и весьма распространенный вариационный метод решения эллиптических уравнений. Опишем его, по-прежнему используя как пример задачу (1.1), (1.2).

Пусть заданы функции  $\varphi_1(x), \varphi_2(x), \dots, \varphi_n(x)$ , такие, что  $\varphi_i(0) = \varphi_i(l) = 0$ ,  $i = 1, 2, \dots, n$ . Эти функции называют координатными или базисными функциями метода Ритца.

Приближенное решение задачи о минимуме функционала будем искать в виде линейной комбинации

$$u_n(x) = \sum_{i=1}^n c_i \varphi_i(x),$$

коэффициенты которой и подлежат определению. А именно, они находятся так, чтобы функция  $n$  вещественных переменных

$$\Phi(c_1, c_2, \dots, c_n) = F\left(\sum_{i=1}^n c_i \varphi_i\right)$$

принимала минимальное значение. Запишем необходимые условия минимума:

$$\frac{\partial F\left(\sum_{i=1}^n c_i \varphi_i\right)}{\partial c_k} = 0, \quad k = 1, 2, \dots, n. \quad (3.1)$$

Система уравнений (3.1) называется системой уравнений метода Ритца.

Видно, что сама формулировка метода Ритца никак не привязана к конкретному виду функционала  $F$ , и, действительно, область применения этого метода чрезвычайно широка.

В рассматриваемом нами случае систему уравнений (3.1) нетрудно записать более подробно. Имеем:

$$F\left(\sum_{i=1}^n c_i \varphi_i\right) = \frac{1}{2} \int_0^l \left[ p \left(\sum_{i=1}^n c_i \varphi_i'\right)^2 + q \left(\sum_{i=1}^n c_i \varphi_i\right)^2 \right] dx - \int_0^l f \sum_{i=1}^n c_i \varphi_i dx.$$

Легко находится, что

$$\frac{\partial F(u_n)}{\partial c_k} = \int_0^l \left[ p \sum_{i=1}^n c_i \varphi_i' \varphi_k' + q \sum_{i=1}^n c_i \varphi_i \varphi_k \right] dx - \int_0^l f \varphi_k dx.$$

Следовательно, система (3.1) принимает вид:

$$\sum_{i=1}^n a_{ki} c_i = b_k, \quad k = 1, 2, \dots, n, \quad (3.2)$$

где

$$a_{ki} = \int_0^l [p \varphi_i' \varphi_k' + q \varphi_i \varphi_k] dx, \quad (3.3)$$

$$b_k = \int_0^l f \varphi_k dx. \quad (3.4)$$

Понятно, что фактическое построение системы требует умения вычислять интегралы (3.3), (3.4). На практике для этого, обычно, применяют квадратурные формулы.

Свойства системы (3.2) и построенного с ее помощью приближенного решения  $u_n$ , конечно, определяются выбором координатной системы  $\varphi_1, \varphi_2, \dots, \varphi_n$ . Если координатная система выбрана удачно, то с увеличением  $n$  точность приближенного метода улучшается. Можно показать, что если координатная система линейно независима, то матрица системы линейных уравнений (3.2) положительно определена. Отсюда, в частности, вытекает однозначная разрешимость системы метода Ритца.

В качестве примеров координатных систем метода Ритца укажем:

а) полиномиальную систему

$$\varphi_0(x) = x(l-x), \quad \varphi_1(x) = x(l-x)x, \quad \varphi_2(x) = x(l-x)x^2, \dots$$

$$\varphi_n(x) = x(l-x)x^n, \dots;$$

б) тригонометрическую систему

$$\varphi_k(x) = \sin \frac{\pi k x}{l}, \quad k = 1, 2, \dots, n, \dots$$

При использовании указанных, или аналогичных, координатных систем матрица метода Ритца оказывается заполненной, т. е. все ее элементы отличны от нуля. Это принципиально отличает метод Ритца от разностного метода, при использовании которого матрица системы линейных уравнений — разреженная матрица. Правда, при удачном выборе координатной системы метода Ритца можно ограничиться небольшими значениями  $n$ , т. е. система будет иметь небольшие размеры и ее решение не вызывает затруднений.

**4. Метод конечных элементов.** В 1943 г. Р. Курантом было замечено, что при специальном выборе базисных функций метод Ритца приводит к системам линейных уравнений, по свойствам весьма близким к разностным уравнениям. В простейших случаях метод совпадает с разностным. Впоследствии (в 50-х годах) метод, предложенный Курантом, был переоткрыт инженерами, существенно развит и обобщен. В настоящее время этот метод, называемый методом конечных элементов (МКЭ), принадлежит к числу наиболее распространенных методов решения эллиптических уравнений и систем, в частности, возникающих при математическом моделировании упругих конструкций. Мы опишем простейший вариант метода конечных элементов на примере задачи (1.1), (1.2).

Построим на отрезке  $[0, l]$  сетку (вообще говоря, неравномерную)

$$\omega = \{x_0 = 0 < x_1 < x_2 < \dots < x_N = l\}.$$

Свяжем с каждой внутренней точкой сетки  $x_k$  непрерывную функцию  $\varphi_k(x)$ , линейную на каждом интервале  $(x_{i-1}, x_i)$ ,  $i = 1, 2, \dots, N$ , и такую, что

$$\varphi_k(x_i) = \begin{cases} 1, & i = k, \\ 0, & i \neq k, \end{cases} \quad i = 0, 1, \dots, N.$$

Легко написать аналитическое выражение для  $\varphi_k(x)$ , но в дальнейшем оно нам не потребуется. Функции  $\varphi_k(x)$  принято называть базисными функциями с локальным носителем, поскольку множество

точек, где  $\varphi_k(x) \neq 0$ , (носитель функции  $\varphi_k(x)$ ) — малый интервал (в рассматриваемом случае — два шага сетки), содержащий точку  $x_k$ . Функции  $\varphi_k(x)$  иногда называют функциями Куранта.

Будем искать решение задачи на минимум функционала  $F$  в виде

$$y(x) = \sum_{k=1}^{N-1} c_k \varphi_k(x).$$

Нетрудно убедиться, что  $y(x)$  — кусочно линейная функция, равная нулю на концах отрезка  $[0, l]$ , причем  $c_k = y(x_k) = y_k$ ,  $k = 1, 2, \dots, N-1$ , то есть можно написать

$$y(x) = \sum_{k=1}^{N-1} y_k \varphi_k(x).$$

Система метода Ритца, по-прежнему, имеет вид (3.2). Элементы  $a_{ki}$  матрицы этой системы, очевидно, отличны от нуля лишь при условии, что носители функций  $\varphi_k(x), \varphi_i(x)$  пересекаются, то есть  $|k-i| < 2$ . Это означает, что матрица системы трехдиагональна, как и в случае разностного метода.

Рассмотрим самый простой частный случай. Пусть  $p(x) \equiv 1$ ,  $q(x) \equiv 0$ , сетка равномерна. Элементы матрицы системы метода Ритца будут при этом равны

$$a_{ki} = \int_0^l \varphi'_k(x), \varphi'_i(x) dx.$$

Фиксируем некоторое  $k$ ,  $2 \leq k \leq N-2$ . Ненулевыми коэффициентами в  $k$ -м уравнении системы метода Ритца будут лишь

$$a_{kk-1} = \int_0^l \varphi'_k(x), \varphi'_{k-1}(x) dx, \quad a_{kk} = \int_0^l \varphi'^2_k(x) dx,$$

$$a_{kk+1} = \int_0^l \varphi'_k(x), \varphi'_{k+1}(x) dx.$$

Подсчитаем их значения. Ясно, что

$$a_{kk-1} = \int_{x_{k-1}}^{x_k} \varphi'_k(x), \varphi'_{k-1}(x) dx,$$



причем

$$\varphi'_k(x) = \frac{1}{h} \text{ при } x \in (x_{k-1}, x_k), \quad \varphi'_{k-1}(x) = -\frac{1}{h} \text{ при } x \in (x_{k-1}, x_k).$$

Отметим, что для вычисления производных нет нужды выписывать аналитические выражения базисных функций. Достаточно заметить, что они линейны на рассматриваемом интервале и, следовательно, их производные совпадают с разностными отношениями, вычисленными, например, по точкам  $x_{k-1}$ ,  $x_k$ . Таким образом, очевидно,  $a_{kk-1} = -1/h$ . Точно так же  $a_{kk+1} = -1/h$ . Далее

$$a_{kk} = \int_{x_{k-1}}^{x_{k+1}} \varphi'^2_k(x) dx = \int_{x_{k-1}}^{x_k} \frac{1}{h^2} dx + \int_{x_k}^{x_{k+1}} \frac{1}{h^2} dx = \frac{2}{h},$$

$$b_k = \int_0^l f \varphi_k(x) dx = \int_{x_{k-1}}^{x_{k+1}} f \varphi_k(x) dx.$$

Вычислим последний интеграл приближенно, полагая  $f(x) \approx f(x_k)$ . Тогда

$$b_k = f(x_k) \int_{x_{k-1}}^{x_{k+1}} \varphi_k(x) dx = hf(x_k),$$

и, следовательно,  $k$ -е уравнение системы метода Рунге принимает вид

$$-\frac{y_{k+1} - 2y_k + y_{k-1}}{h^2} = f_k,$$

то есть совпадает с разностной аппроксимацией уравнения (1.1).

**5. Метод конечных элементов для эллиптических уравнений.** Опишем этот метод на примере задачи Дирихле для уравнения Пуассона:

$$-\Delta u = f, \quad (x, y) \in \Omega, \quad (5.1)$$

$$u(x, y) = 0, \quad (x, y) \in \Gamma. \quad (5.2)$$

Для простоты предполагается, что область  $\Omega$  — единичный квадрат:  $\Omega = [0 < x, y < 1]$ ,  $\Gamma$  — граница  $\Omega$ . Метод конечных элементов (МКЭ) основан на вариационной формулировке задачи (5.1), (5.2). Рассмотрим функционал

$$F(u) = \frac{1}{2} \int_{\Omega} \left[ \left( \frac{\partial u}{\partial x} \right)^2 + \left( \frac{\partial u}{\partial y} \right)^2 \right] dx dy - \int_{\Omega} f u dx dy.$$

ЗАМЕЧАНИЕ 5.1. Если задачу (5.1), (5.2) интерпретировать как задачу о равновесии мембраны, то  $F(u)$  — потенциальная энергия системы мембрана — внешние силы на перемещении  $u$ .

Можно показать, что если  $u$  — решение задачи (5.1), (5.2), то  $u$  доставляет минимальное значение функционалу  $F$  на множестве дважды непрерывных функций, удовлетворяющих граничному условию (5.2), и наоборот, функция, доставляющая минимальное значение функционалу  $F$  на множестве функций, удовлетворяющих граничному условию (5.2), — решение задачи (5.1), (5.2). При этом используются, по существу, точно такие же рассуждения, какие были нами проведены для обыкновенного дифференциального уравнения второго порядка.

Таким образом, можно строить приближенное решение, отправляясь от задачи минимизации функционала  $F$ . Чаще всего для этого используется метод Ритца: приближенное решение разыскивается в виде

$$u_n = \sum_{k=1}^n c_k \varphi_k(x, y),$$

где  $\varphi_k(x, y)$  — некоторые заданные функции, удовлетворяющие граничному условию (5.2) (их называют — базисными), а коэффициенты  $c_k$  подлежат определению. При этом их подбирают так, чтобы доставить функционалу  $F(u_n)$  минимальное значение. Получим подробное выражение для  $F(u_n)$ . Имеем:

$$\begin{aligned} F(u_n) &= \frac{1}{2} \int_{\Omega} \left[ \left( \sum_{k=1}^n c_k \frac{\partial \varphi_k}{\partial x} \right)^2 + \left( \sum_{k=1}^n c_k \frac{\partial \varphi_k}{\partial y} \right)^2 \right] dx dy - \\ &\quad - \int_{\Omega} \sum_{k=1}^n c_k \varphi_k f dx dy = \\ &= \frac{1}{2} \int_{\Omega} \sum_{k,l=1}^n c_k c_l \left( \frac{\partial \varphi_k}{\partial x} \frac{\partial \varphi_l}{\partial x} + \frac{\partial \varphi_k}{\partial y} \frac{\partial \varphi_l}{\partial y} \right) dx dy - \\ &\quad - \int_{\Omega} \sum_{k=1}^n c_k \varphi_k f dx dy = \end{aligned}$$

$$= \frac{1}{2} \sum_{k,l=1}^n a_{kl} c_k c_l - \sum_{k=1}^n c_k f_k \equiv \Phi(c_1, c_2, \dots, c_n),$$

где

$$a_{kl} = \int_{\Omega} \left( \frac{\partial \varphi_k}{\partial x} \frac{\partial \varphi_l}{\partial x} + \frac{\partial \varphi_k}{\partial y} \frac{\partial \varphi_l}{\partial y} \right) dx dy, \quad f_k = \int_{\Omega} \varphi_k f dx dy. \quad (5.3)$$

Запишем необходимые условия минимума функции  $\Phi$ :

$$\frac{\partial \Phi}{\partial c_k} = \sum_{l=1}^n a_{kl} c_l - f_k = 0, \quad k = 1, 2, \dots, n.$$

В результате получаем систему линейных алгебраических уравнений для определения коэффициентов  $c_k$ ,  $k = 1, 2, \dots, n$ :

$$\sum_{l=1}^n a_{kl} c_l = f_k, \quad k = 1, 2, \dots, n. \quad (5.4)$$

Построение системы (5.4) связано с вычислением интегралов, определяющих  $a_{kl}$ ,  $f_k$ . Эта процедура оказывается более или менее трудоемкой в зависимости от выбора функций  $\varphi_k$ .

Классические примеры выбора базисных функций:

$$1) \quad \varphi_{kl}(x, y) = x(1-x)y(1-y)x^k y^l, \quad k, l = 1, 2, \dots, n$$

(первые множители здесь служат для удовлетворения граничным условиям),

$$2) \quad \varphi_{kl}(x, y) = \sin \pi k x \sin \pi l y, \quad k, l = 1, 2, \dots, n.$$

Отметим, что для удобства здесь принята двойная нумерация базисных функций.

Метод конечных элементов основан на использовании специальных базисных функций. При этом матрица системы метода Рунца оказывается разреженной, т. е. большинство ее элементов — нули. Аналогичное имеет место и для существенно более общих эллиптических уравнений, чем уравнение Пуассона.

Построим на области  $\Omega$  квадратную сетку с шагом  $h$ . Каждую ячейку сетки разделим на два треугольника диагональю, параллельной биссектрисе первого координатного угла. Получим разбиение области  $\Omega$  на треугольники, т. е. триангуляцию области.

Рассмотрим произвольную внутреннюю точку сетки  $(x_k, y_l) = (kh, lh)$ . Обозначим через  $\Omega_{kl}$  объединение всех треугольников триангуляции с вершиной в этой точке (их — шесть). Обозначим далее через  $\varphi_{kl}(x, y)$  функцию, равную нулю вне области  $\Omega_{kl}$ , равную единице в точке  $(x_k, y_l)$ , непрерывную на  $\Omega_{kl}$  и линейную на каждом треугольнике триангуляции. Такая функция называется функцией Куранта («шапочкой Куранта»). Используем эти функции при построении приближенного решения по методу Рунца:

$$u^h(x, y) = \sum_{k,l=1}^{N-1} c_{kl} \varphi_{kl}(x, y).$$

Заметим, что

$$u^h(x_i, y_j) = \sum_{k,l=1}^{N-1} c_{kl} \varphi_{kl}(x_i, y_j) = c_{ij},$$

то есть

$$u^h(x, y) = \sum_{k,l=1}^{N-1} u^h(x_k, y_l) \varphi_{kl}(x, y).$$

Это значит, что коэффициентами в разложении приближенного решения служат значения приближенного решения в точках сетки. Кроме того отметим, что приближенное решение — линейная функция на каждом треугольнике триангуляции.

Построим систему уравнений для определения значений  $u^h(x_k, y_l)$ . Для этого нужно только подсчитать коэффициенты и правую часть в системе метода Рунца, а именно,

$$a_{klk'l'} = \int_{\Omega} \left( \frac{\partial \varphi_{kl}}{\partial x} \frac{\partial \varphi_{k'l'}}{\partial x} + \frac{\partial \varphi_{kl}}{\partial y} \frac{\partial \varphi_{k'l'}}{\partial y} \right) dx dy, \quad f_{kl} = \int_{\Omega} \varphi_{kl} f dx dy.$$

Понятно, что  $a_{klk'l'}$  может быть не нулем только в том случае, когда области  $\Omega_{kl}$ ,  $\Omega_{k'l'}$  имеют общие точки. Для заданных  $k, l$  таких областей  $\Omega_{k'l'}$  шесть, а именно:

$$\Omega_{k-1,l}, \Omega_{k-1,l-1}, \Omega_{k,l-1}, \Omega_{k+1,l}, \Omega_{k+1,l+1}, \Omega_{k,l+1}.$$

При вычислении конкретных значений коэффициентов  $a_{klk'l'}$  потребуются значения производных функции  $\varphi_{kl}$ . Пронумеруем треугольники, принадлежащие  $\Omega_{kl}$ . Поскольку функция  $\varphi_{kl}$  линейна на каждом из указанных треугольников, то ее первые производные постоянны

на каждом из этих треугольников, причем производная по  $x$  равна разностному отношению, вычисленному по любым двум точкам, лежащим на прямой, параллельной оси  $x$ ; аналогично вычисляется производная по  $y$ .

Сведем результаты этих очевидных вычислений в таблицу ( $n$  — номер треугольника).

$n$	$\frac{\partial \varphi_{kl}}{\partial x}$	$\frac{\partial \varphi_{kl}}{\partial y}$
1	0	$-1/h$
2	$-1/h$	$1/h$
3	$-1/h$	0
4	0	$-1/h$
5	$1/h$	$-1/h$
6	$1/h$	0

Теперь совершенно очевидно, что

$$\begin{aligned} a_{k,l,k-1,l-1} &= \int_{\Delta_6 \cap \Delta_1} \left( \frac{\partial \varphi_{kl}}{\partial x} \frac{\partial \varphi_{k-1,l-1}}{\partial x} + \frac{\partial \varphi_{kl}}{\partial y} \frac{\partial \varphi_{k-1,l-1}}{\partial y} \right) dx dy = \\ &= \frac{h^2}{2} \left( \frac{1}{h} \cdot 0 + 0 \cdot \frac{-1}{h} \right) + \frac{h^2}{2} \left( 0 \cdot \frac{-1}{h} + \frac{-1}{h} \cdot 0 \right) = 0. \end{aligned}$$

Поясним, что  $h^2/2$  — площадь треугольника  $\Delta_k$ ; значения производных функции  $\varphi_{k-1,l-1}$  на  $\Delta_6$  совпадают со значениями производных  $\varphi_{k,l}$  на  $\Delta_4$ ; значения производных функции  $\varphi_{k-1,l-1}$  на  $\Delta_1$  совпадают со значениями производных  $\varphi_{k,l}$  на  $\Delta_3$ .

Точно так же получаем  $a_{k,l,k+1,l+1} = 0$ ,

$$a_{k,l,k,l-1} = \int_{\Delta_1 \cap \Delta_2} \left( \frac{\partial \varphi_{kl}}{\partial x} \frac{\partial \varphi_{k,l-1}}{\partial x} + \frac{\partial \varphi_{kl}}{\partial y} \frac{\partial \varphi_{k,l-1}}{\partial y} \right) dx dy = -1,$$

$$a_{k,l,k-,l} = -1, \quad a_{k,l,k,l+1} = -1, \quad a_{k,l,k+1,l} = -1,$$

$$a_{k,l,k,l} = \int_{\Omega_{kl}} \left( \frac{\partial \varphi_{kl}}{\partial x} \frac{\partial \varphi_{k,l}}{\partial x} + \frac{\partial \varphi_{kl}}{\partial y} \frac{\partial \varphi_{k,l}}{\partial y} \right) dx dy = 4.$$

Это означает, что уравнение с номером  $k, l$  системы Ритца записывается в виде

$$\frac{4u_{k,l}^h - u_{k-1,l}^h - u_{k+1,l}^h - u_{k,l-1}^h - u_{k,l+1}^h}{h^2} = g_{k,l}^h, \quad (5.5)$$

где

$$g_{k,l}^h = \frac{1}{h^2} \int_{\Omega_{kl}} f(x, y) \varphi_{kl}(x, y) dx dy.$$

Заметим, что при малом  $h$  справедливо приближенное равенство

$$\int_{\Omega_{kl}} f(x, y) \varphi_{kl}(x, y) dx dy \approx f(x_k, y_l) \int_{\Omega_{kl}} \varphi_{kl}(x, y) dx dy = f(x_k, y_l) h^2.$$

Поясним, что для вычисления интеграла от  $\varphi_{kl}(x, y)$  по  $\Omega_{kl}$  достаточно заметить, что он равен объему пирамиды с основанием  $\Omega_{kl}$ , площадь которого равна  $6h^2/2 = 3h^2$ , и высотой, равной единице. Таким образом, уравнение (5.5) приближенно представляется в виде

$$\frac{4u_{k,l}^h - u_{k-1,l}^h - u_{k+1,l}^h - u_{k,l-1}^h - u_{k,l+1}^h}{h^2} = f_{k,l}.$$

Очевидно, точно такое же уравнение получается при замене в уравнении Пуассона (5.1) производных разделенными разностями

$$\frac{\partial^2 u}{\partial x^2} \approx \frac{u(x_{k-1,l}) - 2u(x_{k,l}) + u(x_{k+1,l}))}{h^2},$$

$$\frac{\partial^2 u}{\partial y^2} \approx \frac{u(x_{k,l-1}) - 2u(x_{k,l}) + u(x_{k,l+1}))}{h^2}.$$

Записывая уравнения (5.5) во внутренних точках сетки, т. е. при  $k, l = 1, 2, \dots, N - 1$ , и присоединяя к ним граничные условия, соответствующие (5.2):

$$\begin{aligned} u^h(0, y_l) = 0, \quad u^h(1, y_l) = 0, \quad l = 1, 2, \dots, N - 1, \\ u^h(x_k, 0) = 0, \quad u^h(x_k, 1) = 0, \quad k = 1, 2, \dots, N - 1, \end{aligned} \quad (5.6)$$

получим полную систему линейных алгебраических уравнений для отыскания приближенного решения в точках сетки.

**ЗАМЕЧАНИЕ 5.2.** Совершенно аналогично строится система метода конечных уравнений для произвольных областей: сначала область аппроксимируется многоугольником, затем производят триангуляцию, т. е. разбиение многоугольника на достаточно малые треугольники, далее определяются базисные функции  $\varphi_{kl}$ , при этом, конечно, области  $\Omega_{kl}$  могут иметь достаточно сложное строение.

### 6. Итерационные методы решения сеточных уравнений.

Для решения системы уравнений МКЭ (сеточных уравнений) при умеренных значениях числа неизвестных чаще всего применяют прямые методы, обычно метод Гаусса и его модификации. При очень большом числе неизвестных более эффективными становятся итерационные методы. Опишем простейшие из них, используя в качестве примера систему (5.5).

**6.1.** Метод Якоби (метод простой итерации). Для сокращения записей используем обозначения  $y_{kl} = u^h(x_k, y_l)$ ,  $b_{kl} = g_{kl}^h$ . Зададимся некоторым начальным приближением  $y_{kl}^0$ ,  $k, l = 0, 1, \dots, N$ , удовлетворяющим граничным условиям (5.6). Чаще всего полагают  $y_{kl}^0 = 0$ . Отправляясь от  $y_{kl}^0$ , строим последовательность сеточных функций  $y_{kl}^s$ ,  $s = 1, 2, \dots$ , по рекуррентным формулам:

$$y_{kl}^{s+1} = \frac{1}{4} (y_{k+1,l}^s + y_{k-1,l}^s + y_{k,l+1}^s + y_{k,l-1}^s + h^2 b_{kl}),$$

$$k, l = 1, 2, \dots, N-1, \quad (6.1)$$

$$y_{0l}^{s+1} = 0, \quad y_{Nl}^{s+1} = 0, \quad l = 1, 2, \dots, N-1,$$

$$y_{k0}^{s+1} = 0, \quad y_{kN}^{s+1} = 0, \quad k = 1, 2, \dots, N-1,$$

$s = 0, 1, 2, \dots$

В ходе вычислений контролируется невязка

$$\frac{4y_{k,l}^s - y_{k-1,l}^s - y_{k+1,l}^s - y_{k,l-1}^s - y_{k,l+1}^s}{h^2} - b_{k,l}.$$

Вычисления останавливают, если при некотором  $s$  значения невязки при всех  $k, l = 1, 2, \dots, N-1$  по модулю не превышают заданной величины  $\varepsilon$ , определяющей точность вычислений.

**6.2.** Метод Зейделя. При вычислениях по формуле (6.1) порядок перебора точек сетки, по крайней мере теоретически, может быть произвольным. Условимся теперь перебирать точки в порядке возрастания индексов  $k, l$  и модифицируем формулы (6.1) следующим образом:

$$y_{kl}^{s+1} = \frac{1}{4} (y_{k+1,l}^s + y_{k-1,l}^{s+1} + y_{k,l+1}^s + y_{k,l-1}^{s+1} + h^2 b_{kl}),$$

$$k, l = 1, 2, \dots, N-1, \quad (6.2)$$

$$y_{0l}^{s+1} = 0, \quad y_{Nl}^{s+1} = 0, \quad l = 1, 2, \dots, N-1,$$

$$y_{k0}^{s+1} = 0, \quad y_{kN}^{s+1} = 0, \quad k = 1, 2, \dots, N-1,$$

$s = 0, 1, 2, \dots$  Этот метод сходится не медленнее, чем метод Якоби (на практике, как правило, быстрее).

**6.3.** Метод верхней релаксации. Значительного ускорения сходимости метода Зейделя можно добиться, несколько модифицируя формулы (6.2) введением итерационного параметра:

$$y_{kl}^{s+1} = (1 - \omega)y_{kl}^s + \omega \frac{1}{4} \left( y_{k+1,l}^s + y_{k-1,l}^{s+1} + y_{k,l+1}^s + y_{k,l-1}^{s+1} + h^2 b_{kl} \right), \quad (6.3)$$

$$k, l = 1, 2, \dots, N - 1,$$

$$\begin{aligned} y_{0l}^{s+1} &= 0, & y_{Nl}^{s+1} &= 0, & l &= 1, 2, \dots, N - 1, \\ y_{k0}^{s+1} &= 0, & y_{kN}^{s+1} &= 0, & k &= 1, 2, \dots, N - 1, \end{aligned}$$

$s = 0, 1, 2, \dots$  Понятно, что при  $\omega = 1$  метод (6.3) превращается в метод Зейделя. Оптимальное значение параметра  $\omega$  зависит от сетки. На практике, обычно,  $\omega$  полагают близким к 1,8.

**7. Разностные методы решения нестационарных задач математической физики.** Рассмотрим первую краевую задачу для уравнения теплопроводности стержня:

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left( p \frac{\partial u}{\partial x} \right) - qu + f(x, t), \quad 0 < x < l, \quad t > 0, \quad (7.1)$$

$$u(x, 0) = u_0(x), \quad 0 < x < l, \quad (7.2)$$

$$u(0, t) = u(l, t) = 0, \quad t > 0. \quad (7.3)$$

**7.1.** Явная схема. Построим на области  $\Omega = [0 < x < l, t > 0]$  сетку  $\omega$  с шагами  $h$  по оси  $x$ ,  $\tau$  — по оси  $t$ . Значение сеточной функции  $y = y(x_i, t_j)$  в точке  $(x_i, t_j) = (ih, t_j)$ ,  $0 \leq i \leq N, Nh = l, j \geq 0$  обозначим через  $y_i^j$ . Как и ранее, пусть

$$y_{x,i}^j = \frac{y(x_{i+1}, t_j) - y(x_i, t_j)}{h} \quad \text{— разность вперед по } x,$$

$$y_{\bar{x},i}^j = \frac{y(x_i, t_j) - y(x_{i-1}, t_j)}{h} \quad \text{— разность назад по } x.$$

Положим еще

$$y_{t,i}^j = \frac{y(x_i, t_{j+1}) - y(x_i, t_j)}{h} \quad \text{— разность вперед по } t.$$

Действуя теперь по аналогии с обыкновенным дифференциальным уравнением, заменим в точках сетки  $(x_i, t_j)$  уравнение (7.1) — разностным

$$y_{t,i}^j = (ay_{\bar{x}}^j)_{x_i} - by_i^j + \varphi_i^j, \quad 1 \leq i \leq N - 1, \quad j \geq 0. \quad (7.4)$$



В результате, получим систему линейных алгебраических уравнений. Эта система, однако, неполная. Нужно еще учесть начальные и граничные условия:

$$y(x_i, 0) = u_0(x_i), \quad 0 \leq i \leq N, \quad (7.5)$$

$$y(0, t_j) = y(x_N, t_j) = 0, \quad j \geq 0. \quad (7.6)$$

Совокупность линейных алгебраических уравнений (7.4)–(7.6) называется разностной схемой для задачи (7.1)–(7.3). Это — явная схема. Вычисления по указанной схеме проводятся так. Значения  $y_i^0$  при всех  $i$ ,  $0 \leq i \leq N$ , известны, так как задано начальное условие (7.5). Предположим, что уже найдены значения  $y_i^j$ ,  $0 \leq i \leq N$ . Совокупность точек сетки  $(x_i, t_j)$  при фиксированном  $j$  называется  $j$ -м временным слоем. Напишем уравнение (7.4) во всех внутренних точках  $j$ -го временного слоя чуть более подробно:

$$\frac{y_i^{j+1} - y_i^j}{\tau} = (ay_{\bar{x}}^j)_{x_i} - by_i^j + \varphi_i^j, \quad 1 \leq i \leq N - 1,$$

или

$$y_i^{j+1} = y_i^j + \tau \left[ (ay_{\bar{x}}^j)_{x_i} - by_i^j + \varphi_i^j \right], \quad 1 \leq i \leq N - 1. \quad (7.7)$$

В правой части последнего уравнения только известные величины. Поэтому значения  $y$  во всех точках  $j+1$  временного слоя вычисляются по явным формулам.

**7.2.** Неявная схема. По аналогии с (7.4) можно написать разностное уравнение

$$y_{t,i}^j = (ay_{\bar{x}}^{j+1})_{x_i} - by_i^{j+1} + \varphi_i^j, \quad 1 \leq i \leq N - 1, \quad j \geq 0. \quad (7.8)$$

Присоединяя начальные и граничные условия

$$y(x_i, 0) = u_0(x_i), \quad 0 \leq i \leq N, \quad (7.9)$$

$$y(0, t_j) = y(x_N, t_j) = 0, \quad j \geq 0, \quad (7.10)$$

вновь получим полную систему линейных алгебраических уравнений. Эта система называется чисто неявной (чаще, просто, неявной) разностной схемой для уравнения теплопроводности. Построение решения этой системы кардинально отличается от случая явной схемы. Запишем разностное уравнение (7.8) подробнее в точках  $j$ -го временного слоя:

$$\frac{y_i^{j+1} - y_i^j}{\tau} = (ay_{\bar{x}}^{j+1})_{x_i} - by_i^{j+1} + \varphi_i^j, \quad 1 \leq i \leq N - 1,$$

или

$$y_i^{j+1} - \tau(ay_x^{j+1})_{x_i} + \tau by_i^{j+1} = y_i^j + \tau\varphi_i^j, \quad 1 \leq i \leq N-1.$$

В силу граничных условий

$$y_0^{j+1} = y_N^{j+1} = 0.$$

Более подробная запись полученной системы уравнений такова:

$$A_i y_{i+1}^{j+1} - B_i y_i^{j+1} + C_i y_{i-1}^{j+1} = -F_i^j, \quad (7.11)$$

$$y_0^{j+1} = y_N^{j+1} = 0,$$

где

$$A_i = \frac{\tau a_{i+1}}{h^2}, \quad C_i = \frac{\tau a_i}{h^2}, \quad B_i = 1 + A_i + C_i + b_i, \quad F_i^j = (y_i^j + \tau\varphi_i^j). \quad (7.12)$$

Систему (7.11) можно решать методом прогонки. Условие реализуемости метода прогонки выполнено, так как, очевидно,

$$A_i + C_i < B_i, \quad 1 \leq i \leq N-1.$$

**7.3.** Схемы с весами. По аналогии с явной и неявной схемами можно написать целое семейство схем (это так называемые схемы с весами):

$$y_{t,i}^j = (ay_x^{(\sigma)})_{x_i} - by_i^{(\sigma)} + \varphi_i^j, \quad 1 \leq i \leq N-1, \quad j \geq 0. \quad (7.13)$$

$$y(x_i, 0) = u_0(x_i), \quad 0 \leq i \leq N, \quad (7.14)$$

$$y(0, t_j) = y(x_N, t_j) = 0, \quad j \geq 0, \quad (7.15)$$

где  $y_i^{(\sigma)} = \sigma y_i^{j+1} + (1 - \sigma)y_i^j$ ,  $\sigma$  — число, называемое весом слоя.

При  $\sigma = 0$  как частный случай мы получаем явную схему, а при  $\sigma = 1$  — чисто неявную схему. Все схемы при  $\sigma \neq 0$  неявные. Их решения строятся при помощи метода прогонки. Среди схем с весами особо выделяется схема с  $\sigma = 1/2$ . Эта схема называется симметричной схемой с весами, или схемой Кранка — Никольсон. Если провести вполне уместную аналогию с обыкновенными дифференциальными уравнениями, то явная схема аналогична явному методу Эйлера, неявная схема — неявному методу Эйлера, а схема Кранка — Никольсон — методу трапеций.

Погрешность аппроксимации схемы с весами. Основной характеристикой схемы является погрешность аппроксимации, т. е. сеточная функция

$$\psi_i^j = u_{t,i}^j - (au_x^{(\sigma)})_{x_i} + bu^{(\sigma)} - \varphi_i^j,$$

где  $u$  — точное решение задачи (7.1)–(7.3). Понятно, что она зависит от  $\sigma$  и способа выбора функций  $a$ ,  $b$ ,  $\varphi$ . Покажем, что если  $a$ ,  $b$  выбраны в соответствии с (1.4), а  $\varphi_i^j = f(x_i, t_j + \tau/2)$ , то

$$\psi_i^j = \begin{cases} O(h^2 + \tau), & \sigma \neq 1/2, \\ O(h^2 + \tau^2), & \sigma = 1/2, \end{cases} \quad (7.16)$$

т. е. среди всех схем с весами наилучшую оценку погрешности аппроксимации имеет схема Кранка — Никольсона.

Используя тождество

$$\begin{aligned} u^{(\sigma)} &= \sigma u_i^{j+1} + (1 - \sigma)u_i^j = \frac{u_i^{j+1} + u_i^j}{2} - \left(\sigma - \frac{1}{2}\right) \tau \frac{u_i^{j+1} - u_i^j}{\tau} = \\ &= \frac{u_i^{j+1} + u_i^j}{2} - \left(\sigma - \frac{1}{2}\right) \tau u_{t,i}^j. \end{aligned}$$

представим погрешность аппроксимации в виде  $\psi_i^j = \psi_i^{j(0)} + \psi_i^{j(1)}$ , где

$$\begin{aligned} \psi_i^{j(0)} &= u_{t,i}^j - (au_{\bar{x}}^{(1/2)})_{x_i} + bu_i^{(1/2)} - \varphi_i^j, \\ \psi_i^{j(1)} &= \tau \left(\sigma - \frac{1}{2}\right) [(au_{t\bar{x}})_x - bu_t]. \end{aligned}$$

Очевидно,

$$\psi_i^{j(1)} = \begin{cases} O(\tau), & \sigma \neq 1/2, \\ 0, & \sigma = 1/2. \end{cases}$$

Оценим  $\psi_i^{j(0)}$ . По формуле Тейлора имеем:

$$u_i^{j+1} = u_i^{j+1/2} + \frac{\tau}{2} \left(\frac{\partial u}{\partial t}\right)_i^{j+1/2} + \frac{1}{2} \left(\frac{\tau}{2}\right)^2 \left(\frac{\partial^2 u}{\partial t^2}\right)_i^{j+1/2} + \dots, \quad (7.17)$$

$$u_i^j = u_i^{j+1/2} - \frac{\tau}{2} \left(\frac{\partial u}{\partial t}\right)_i^{j+1/2} + \frac{1}{2} \left(\frac{\tau}{2}\right)^2 \left(\frac{\partial^2 u}{\partial t^2}\right)_i^{j+1/2} - \dots, \quad (7.18)$$

то есть

$$u^{(1/2)} = u^{j+1/2} + O(\tau^2).$$

Используя еще раз разложения (7.17)–(7.18), нетрудно проверить, что

$$u_t = \frac{u^{j+1} - u^j}{\tau} = \left(\frac{\partial u}{\partial t}\right)_i^{j+1/2} + O(\tau^2). \quad (7.19)$$

Повторяя теперь выкладки, проведенные на с. 89, получим:

$$\begin{aligned} -(au_{\bar{x}}^{(1/2)})_{x_i} + bu_i^{(1/2)} &= -(au_{\bar{x}}^{j+1/2})_{x_i} + bu_i^{j+1/2} + O(\tau^2) = \\ &= -(pu')'(t_{j+1/2}, x_i) + bu'(t_{j+1/2}, x_i) + O(h^2) + O(\tau^2), \end{aligned}$$

откуда, очевидно, вытекает, что  $\psi_i^{j(0)} = O(h^2 + \tau^2)$ . Таким образом, оценка (7.16) доказана.

Малость погрешности аппроксимации еще не гарантирует близости точного и приближенного решения. Нужно еще, чтобы разностная схема была устойчива. Исследование устойчивости разностных схем для уравнения теплопроводности проведем лишь в частных случаях. Будем предполагать, что  $q(x) \equiv 0$  (это влечет выполнения условия  $b(x) \equiv 0$ ), и рассмотрим только явную ( $\sigma = 0$ ) и чисто неявную ( $\sigma = 1$ ) разностные схемы.

1. Исследование устойчивости явной разностной схемы. Воспользуемся равенством (7.7). Запишем его более подробно:

$$y_i^{j+1} = \left(1 - \frac{\tau}{h^2} (a_{i+1} + a_i)\right) y_i^j + \frac{\tau}{h^2} a_{i+1} y_{i+1}^j + \frac{\tau}{h^2} a_i y_{i-1}^j + \tau \varphi_i^j,$$

$i = 1, 2, \dots, N - 1$ , откуда

$$\begin{aligned} |y_i^{j+1}| &\leq \left|1 - \frac{\tau}{h^2} (a_{i+1} + a_i)\right| |y_i^j| + \frac{\tau}{h^2} a_{i+1} |y_{i+1}^j| + \\ &\quad + \frac{\tau}{h^2} a_i |y_{i-1}^j| + \tau |\varphi_i^j|, \quad i = 1, 2, \dots, N - 1. \end{aligned} \quad (7.20)$$

Предположим, что

$$1 - \frac{\tau}{h^2} (a_{i+1} + a_i) \geq 0, \quad i = 1, 2, \dots, N - 1.$$

Это условие, очевидно, будет выполнено, если

$$\tau \leq \frac{h^2}{2 \max_{0 \leq x \leq l} p(x)}. \quad (7.21)$$

В этом случае знак модуля в первом множителе первого слагаемого справа в (7.20) можно убрать, и, следовательно,

$$|y_i^{j+1}| \leq \max_{0 \leq i \leq N} |y_i^j| + \tau \max_{1 \leq i \leq N-1} |\varphi_i^{j+1}|. \quad (7.22)$$

Выражение в правой части (7.22) не зависит от  $i$ , поэтому

$$\max_{0 \leq i \leq N} |y_i^{j+1}| \leq \max_{0 \leq i \leq N} |y_i^j| + \tau \max_{1 \leq i \leq N-1} |\varphi_i^{j+1}|. \quad (7.23)$$

Неравенство (7.23) выполнено для любого  $j \geq 0$ , следовательно,

$$\max_{0 \leq i \leq N} |y_i^j| \leq \max_{0 \leq i \leq N} |y_i^0| + \tau \sum_{k=0}^{j-1} \max_{1 \leq i \leq N-1} |\varphi_i^k|.$$

Усиливая последнее неравенство, можно написать:

$$\max_{0 \leq i \leq N} |y_i^j| \leq \max_{0 \leq i \leq N} |y_i^0| + t_{j-1} \max_{0 \leq k \leq j-1} \max_{1 \leq i \leq N-1} |\varphi_i^k|. \quad (7.24)$$

Неравенство (7.24) означает устойчивость явной разностной схемы (на конечном отрезке времени). Отсюда сразу вытекает сходимость разностной схемы. Действительно, если  $z = y - u$  — погрешность разностной схемы, то, подставляя  $y = z + u$  в уравнения (7.4)–(7.6), получим разностную схему для погрешности.

$$z_t = (az_{\bar{x}})_x + \psi,$$

$$z_i^0 = 0,$$

$$z_0^j = z_N^j = 0,$$

откуда при выполнении условия (7.21) сразу получим

$$\max_{0 \leq i \leq N} |z_i^j| \leq t_{j-1} \max_{0 \leq k \leq j-1} \max_{1 \leq i \leq N-1} |\psi_i^k| = O(\tau + h^2).$$

2. Исследование устойчивости неявной схемы. Обратимся к уравнению (7.11). Вследствие конечности числа точек сетки на временном слое найдется такой номер  $i_0$ ,  $1 \leq i_0 \leq N - 1$ , что

$$|y_{i_0}^{j+1}| = \max_{0 \leq i \leq N} |y_i^{j+1}|. \quad (7.25)$$

Запишем уравнение (7.11) при  $i = i_0$ . При этом слева оставим только слагаемое, содержащее  $y_{i_0}^{j+1}$ , а остальные — перенесем в левую часть:

$$B_{i_0} y_{i_0}^{j+1} = A_{i_0} y_{i_0+1}^{j+1} + C_{i_0} y_{i_0-1}^{j+1} + \left( y_{i_0}^j + \tau \varphi_{i_0}^j \right),$$

откуда

$$B_{i_0} |y_{i_0}^{j+1}| \leq A_{i_0} |y_{i_0+1}^{j+1}| + C_{i_0} |y_{i_0-1}^{j+1}| + |y_{i_0}^j| + \tau |\varphi_{i_0}^j|,$$

следовательно (см. (7.25)),

$$B_{i_0} \max_{0 \leq i \leq N} |y_i^{j+1}| \leq (A_{i_0} + C_{i_0}) \max_{0 \leq i \leq N} |y_i^{j+1}| + \max_{0 \leq i \leq N} |y_i^j| + \tau \max_{0 \leq i \leq N} |\varphi_i^j|.$$

Используя теперь равенство  $B_{i_0} - A_{i_0} - C_{i_0} = 1 + b_{i_0}$  (см. (7.12)) и условие  $b_i \geq 0$ , получим

$$\max_{0 \leq i \leq N} |y_i^{j+1}| \leq \max_{0 \leq i \leq N} |y_i^j| + \tau \max_{0 \leq i \leq N} |\varphi_i^j|,$$

откуда точно так же, как и в случае явной разностной схемы вытекает неравенство устойчивости неявной разностной схемы и ее сходимости со скоростью  $O(\tau + h^2)$ .

1. При отсутствии внешних источников тепла ( $f(x, t) \equiv 0$ ) неравенство (7.3) принимает вид

$$\max_{0 \leq i \leq N} |y_i^{j+1}| \leq \max_{0 \leq i \leq N} |y_i^j|.$$

Его можно трактовать как принцип максимума для разностной схемы: температура стержня при отсутствии источников тепла не возрастает с ростом времени.

2. Устойчивость явной разностной схемы была нами доказана лишь при условии, что шаги сетки удовлетворяют неравенству (7.21). Это условие на шаги сетки, как можно проверить, является необходимым. При его нарушении решение явной разностной схемы быстро «разбалтывается». Условие (7.21) на практике оказывается особенно обременительным, если коэффициент  $p$  быстро меняется. Тогда шаг по времени приходится брать слишком маленьким, возможно существенно меньшим, чем это диктуется соображениями точности. Это ведет к неоправданному увеличению вычислительной работы. Неявная схема устойчива, как мы показали, при любых шагах сетки (абсолютно устойчива). При ее использовании шаги сетки можно выбирать лишь из соображений точности.

3. Применяемая нами методика (основанная на дискретном принципе максимума) не может быть использована для анализа устойчивости схем с весами при  $\sigma \neq 0$ ,  $\sigma \neq 1$ . Другими, более сложными методами можно показать, что все схемы при  $\sigma \geq 1/2$  абсолютно устойчивы. В частности, абсолютно устойчива и схема Кранка — Никольсона. Все схемы при  $\sigma < 1/2$  условно устойчивы, то есть устойчивы, если шаг по времени достаточно мал по сравнению с  $h^2$ .



- таблицы и графики значений  $y_i$  и  $y_i^k$  с указанием числа итераций, потребовавшихся для достижения заданной точности;
- определение оптимального параметра  $\omega$  для метода релаксации (графики зависимости числа итераций от  $\omega$ );
- графики убывания погрешностей в зависимости от числа итераций и применяемого метода на одном рисунке;
- листинг программы.

### 3. Варианты заданий.

1.  $\alpha = 1, \beta = 1, \gamma = 1.$
2.  $\alpha = 2, \beta = 1, \gamma = 1.$
3.  $\alpha = 1, \beta = 2, \gamma = 1.$
4.  $\alpha = 1, \beta = 1, \gamma = 2.$
5.  $\alpha = 2, \beta = 2, \gamma = 2.$
6.  $\alpha = 3, \beta = 1, \gamma = 1.$
7.  $\alpha = 3, \beta = 2, \gamma = 1.$
8.  $\alpha = 3, \beta = 1, \gamma = 2.$
9.  $\alpha = 3, \beta = 2, \gamma = 2.$
10.  $\alpha = 4, \beta = 1, \gamma = 1.$

## §2. Нелинейные уравнения

На отрезке  $[a, b]$  задана функция вида

$$f(x) = \sum_{n=0}^{\infty} a_n(x).$$



### 1. Построить таблицу обратной к $f(x)$ функции, $F = f^{-1}$

$F_0$	$F_1$	$F_2$	$\dots$	$F_n$
$z_0$	$z_1$	$z_2$	$\dots$	$z_n$

решая уравнения

$$f(z) = F_i, \quad F_i = f(x_0) + i \cdot \frac{f(x_n) - f(x_0)}{n}, \quad i = 0, 1, \dots, n,$$

здесь  $x_i = a + i \cdot h$ ,  $h = \frac{b - a}{n}$ .

Нелинейные уравнения решить итерационными методами:

1. касательных,
2. хорд,
3. секущих.

Во всех итерационных методах в качестве начального приближения к точке  $z_i$  взять  $x_i$  и вычисления продолжать до выполнения условия

$$|r^k| \leq \varepsilon,$$

$r$  — невязка,  $\varepsilon$  — заданное число.

### 2. Отчет должен содержать:

- постановку задачи и исходные данные;
- описание методов решения;
- графики функций  $f(x)$  и  $f^{-1}(x)$  на одном рисунке;
- листинг программы.

### 3. Варианты заданий.

1.  $\sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{(n!)^2}$ ,  $a = 0$ ,  $b = 3$ .
2.  $\frac{2}{\sqrt{\pi}} \sum_{n=0}^{\infty} (-1)^n \frac{x^{2n+1}}{(n!)(2n+1)}$ ,  $a = 0$ ,  $b = 2$ .
3.  $\sum_{n=0}^{\infty} (-1)^n \frac{x^{2n+1}}{(2n+1)(2n+1)!}$ ,  $a = 0$ ,  $b = 4$ .

4.  $\sum_{n=1}^{\infty} (-1)^n \frac{x^{2n}}{2n(2n)!}, \quad a = 0.4, \quad b = 4.$
5.  $\sum_{n=0}^{\infty} (-1)^n \frac{(\pi/2)^{2n} x^{4n+1}}{(2n)!(4n+1)}, \quad a = 0, \quad b = 1.5.$
6.  $\sum_{n=0}^{\infty} (-1)^n \frac{(\pi/2)^{2n+1} x^{4n+3}}{(2n+1)!(4n+1)}, \quad a = 0, \quad b = 1.2.$
7.  $\sum_{n=1}^{\infty} (-1)^n \frac{(x-1)^n}{n^2}, \quad a = 0, \quad b = 2.$
8.  $\frac{x}{2} \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+1)!}, \quad a = 0, \quad b = 3.$
9.  $\left(\frac{x}{2}\right)^2 \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+2)!}, \quad a = 0, \quad b = 4.$
10.  $\left(\frac{x}{2}\right)^3 \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+3)!}, \quad a = 2, \quad b = 6.$
11.  $\left(\frac{x}{2}\right)^4 \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+4)!}, \quad a = 3, \quad b = 7.$

**4. Замечание.** При вычислении ряда  $\sum_{n=0}^{\infty} a_n(x)$  учесть, что каждый последующий член ряда  $a_{n+1}$  получается из предыдущего члена  $a_n$  умножением на величину  $q_n$ , то есть  $a_{n+1} = a_n q_n$ . Это позволит избежать переполнения при вычислении факториалов.

### § 3. Интерполирование функций

На отрезке  $[a, b]$  задана функция вида

$$f(x) = \sum_{n=0}^{\infty} a_n(x).$$

1. Вычислить значения данной функции и ее производной с помощью интерполяционного полинома Лагранжа  $L_n(x)$ . В качестве узлов интерполяции взять:

- 1) равномерно распределенные точки на отрезке  $[a, b]$ ,
- 2) чебышевский набор узлов на отрезке  $[a, b]$ .

При табулировании функции вычислять ряд с точностью до  $10^{-6}$ .

2. Вычислить погрешность интерполирования

$$\varepsilon_1(x) = |f(x) - L_n(x)|, \quad \varepsilon_{1n} = \max_{x \in (a,b)} \varepsilon_1(x),$$

$$\varepsilon_2(x) = |f'(x) - L'_n(x)|, \quad \varepsilon_{2n} = \max_{x \in (a,b)} \varepsilon_2(x).$$

3. Исследовать зависимость погрешности  $\varepsilon_{in}$  от числа узлов интерполяции.

4. Отчет должен содержать:

- постановку задачи и исходные данные;
- описание методов решения;
- графики функций  $f(x)$ ,  $L_n(x)$ ,  $f'(x)$ ,  $L'_n(x)$ ,  $\varepsilon_{in}(x)$ ;
- графики зависимости  $\varepsilon_{in}$  от числа узлов интерполяции;
- листинг программы.

5. Варианты заданий.

1.  $\sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{(n!)^2}$ ,  $a = 0$ ,  $b = 3$ .

2.  $\frac{2}{\sqrt{\pi}} \sum_{n=0}^{\infty} (-1)^n \frac{x^{2n+1}}{(n!)(2n+1)}$ ,  $a = 0$ ,  $b = 2$ .

$$3. \sum_{n=0}^{\infty} (-1)^n \frac{x^{2n+1}}{(2n+1)(2n+1)!}, \quad a=0, \quad b=4.$$

$$4. \sum_{n=1}^{\infty} (-1)^n \frac{x^{2n}}{2n(2n)!}, \quad a=0.4, \quad b=4.$$

$$5. \sum_{n=0}^{\infty} (-1)^n \frac{(\pi/2)^{2n} x^{4n+1}}{(2n)!(4n+1)}, \quad a=0, \quad b=1.5.$$

$$6. \sum_{n=0}^{\infty} (-1)^n \frac{(\pi/2)^{2n+1} x^{4n+3}}{(2n+1)!(4n+1)}, \quad a=0, \quad b=1.2.$$

$$7. \sum_{n=1}^{\infty} (-1)^n \frac{(x-1)^n}{n^2}, \quad a=0, \quad b=2.$$

$$8. \frac{x}{2} \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+1)!}, \quad a=0, \quad b=3.$$

$$9. \left(\frac{x}{2}\right)^2 \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+2)!}, \quad a=0, \quad b=4.$$

$$10. \left(\frac{x}{2}\right)^3 \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+3)!}, \quad a=2, \quad b=6.$$

$$11. \left(\frac{x}{2}\right)^4 \sum_{n=0}^{\infty} (-1)^n \frac{(x^2/4)^n}{n!(n+4)!}, \quad a=3, \quad b=7.$$

**6. Замечание.** При вычислении ряда  $\sum_{n=0}^{\infty} a_n(x)$  учесть, что каждый последующий член ряда  $a_{n+1}$  получается из предыдущего члена  $a_n$  умножением на величину  $q_n$ , то есть  $a_{n+1} = a_n q_n$ . Это позволит избежать переполнения при вычислении факториалов.

## § 4. Численное интегрирование

Вычислить значения интегралов вида

$$I(x) = \int_a^x f(t)dt, \quad I(x) = \int_a^b f(x, t)dt$$

в точках  $x_i = a + i \cdot h$ ,  $i = 0, 1, \dots, n$ , где  $h = (b - a)/n$ , используя составные квадратурные формулы:

1. левых прямоугольников,
2. центральных прямоугольников,
3. трапеции,
5. Симпсона,
6. Гаусса с двумя узлами.

Интеграл вычислить с точностью  $\varepsilon = 10^{-6}$ . Точность вычисления интеграла определяется сравнением результатов при различном числе разбиения отрезка интегрирования. Именно, точность  $\varepsilon$  считается достигнутой, если

$$|S^N(f) - S^{2N}(f)| \leq \varepsilon,$$

здесь  $S^N$  — значение составной квадратурной формулы при разбиении отрезка интегрирования на  $N$  частей.

### 1. Отчет должен содержать:

- постановку задачи и исходные данные;
- описание методов решения и расчетные формулы;
- таблицы значений интегралов с указанием числа разбиений, потребовавшихся для достижения заданной точности;
- листинг программы.

### 2. Варианты заданий.

1.  $\frac{1}{\pi} \int_0^{\pi} \cos(x \cos t) dt$ ,  $a = 0$ ,  $b = 3$ .

$$2. \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt, a = 0, b = 2.$$

$$3. \int_0^x \frac{\sin t}{t} dt, a = 0, b = 4.$$

$$4. \int_0^x \frac{\cos t - 1}{t} dt, a = 0.4, b = 4.$$

$$5. \int_0^x \cos\left(\frac{\pi t^2}{2}\right) dt, a = 0, b = 1.5.$$

$$6. \int_0^x \sin\left(\frac{\pi t^2}{2}\right) dt, a = 0, b = 1.2.$$

$$7. - \int_0^x \frac{\ln t}{1-t} dt, a = 0, b = 2.$$

$$8. \frac{1}{\pi} \int_0^{\pi} \cos(x \sin t - t) dt, a = 0, b = 3.$$

$$9. \frac{1}{\pi} \int_0^{\pi} \cos(x \sin t - 2t) dt, a = 0, b = 4.$$

$$10. \frac{1}{\pi} \int_0^{\pi} \cos(x \sin t - 3t) dt, a = 2, b = 6.$$

$$11. \frac{1}{\pi} \int_0^{\pi} \cos(x \sin t - 4t) dt, a = 3, b = 7.$$

## § 5. Задача Коши для системы обыкновенных дифференциальных уравнений

Решить задачу Коши для системы из двух дифференциальных уравнений первого порядка вида

$$y' = f(t, y), \quad y(0) = y_0, \quad y(t) \in R^2,$$

на отрезке  $[a, b]$ , используя методы Рунге-Кутты с постоянным шагом  $h$ :

1. 2-го порядка точности

$$\begin{aligned} k_1 &= f(t_n, y_n), \\ k_2 &= f(t_n + h, y_n + hk_1), \\ y_{n+1} &= y_n + h(k_1 + k_2)/2; \end{aligned}$$

2. 3-го порядка точности

$$\begin{aligned} k_1 &= f(t_n, y_n), \\ k_2 &= f(t_n + h/3, y_n + h/3k_1), \\ k_3 &= f(t_n + 2/3h, y_n + 2/3hk_2), \\ y_{n+1} &= y_n + h(k_1 + 3k_3)/4; \end{aligned}$$

3. 4-го порядка точности

$$\begin{aligned} k_1 &= f(t_n, y_n), \\ k_2 &= f(t_n + h/4, y_n + h/4k_1), \\ k_3 &= f(t_n + h/2, y_n + h/2k_2), \\ k_4 &= f(t_n + h, y_n + hk_1 - 2hk_2 + 2hk_3), \\ y_{n+1} &= y_n + h(k_1 + 4k_3 + k_4)/6; \end{aligned}$$

4. 5-го порядка точности

$$\begin{aligned} k_1 &= f(t_n, y_n), \\ k_2 &= f(t_n + h/3, y_n + hk_1/3), \\ k_3 &= f(t_n + h/3, y_n + hk_1/6 + hk_2/6), \\ k_4 &= f(t_n + h/2, y_n + h/8k_1 + 3/8hk_3), \\ k_5 &= f(t_n + h, y_n + h/2k_1 - 3/2hk_3 + 2hk_4), \\ y_{n+1} &= y_n + h(k_1 + 4k_4 + k_5)/6. \end{aligned}$$

**1. Отчет должен содержать:**

- постановку задачи и исходные данные;

- описание методов решения;
- графики зависимости максимальной погрешности решения  $\epsilon$  и  $\epsilon/h^k$  от выбранного шага  $h$  для каждого метода ( $k$  – порядок точности метода);
- сделать сравнительный анализ методов.

## 2. Варианты заданий.

1.

$$\begin{aligned} y_1' &= -y_2 + y_1(y_1^2 + y_2^2 - 1), \\ y_2' &= y_1 + y_2(y_1^2 + y_2^2 - 1), \end{aligned}$$

$a = 0, b = 5$ , точное решение:

$$y_1 = \cos(t)/(1 + e^{2t})^{1/2}, \quad y_2 = \sin(t)/(1 + e^{2t})^{1/2}.$$

2.

$$\begin{aligned} y_1' &= -\alpha y_1 - \beta y_2 + (\alpha + \beta - 1)e^{-t}, \\ y_2' &= \beta y_1 - \alpha y_2 + (\alpha + \beta - 1)e^{-t}, \end{aligned}$$

$a = 0, b = 4$ , точное решение:

$$y_1 = y_2 = e^{-t}, \quad \alpha = 2, \quad \beta = 3.$$

3.

$$y_1' = y_2, \quad y_2' = 2y_1^2(1 - 4t^2y_1),$$

$a = 0, b = 5$ , точное решение: (проверьте!)

$$y_1 = 1/(1 + t^2), \quad y_2 = -2t/(1 + t^2)^2.$$

4.

$$\begin{aligned} y_1' &= -\sin(t)/(1 + e^{2t})^{1/2} + y_1(y_1^2 + y_2^2 - 1), \\ y_2' &= \cos(t)/(1 + e^{2t})^{1/2} + y_2(y_1^2 + y_2^2 - 1), \end{aligned}$$

$a = 0, b = 5$ , точное решение:

$$y_1 = \cos(t)/(1 + e^{2t})^{1/2}, \quad y_2 = \sin(t)/(1 + e^{2t})^{1/2}.$$

5.

$$\begin{aligned} y_1' &= -\sin(t)/(1 + e^{2t})^{1/2} + y_1(y_1^2 + y_2^2 - 1), \\ y_2' &= \cos(t)/(1 + e^{2t})^{1/2} + y_2(y_1^2 + y_2^2 - 1), \end{aligned}$$

$a = 0, b = 5$ , точное решение:

$$y_1 = \cos(t)/(1 + e^{2t})^{1/2}, \quad y_2 = \sin(t)/(1 + e^{2t})^{1/2}.$$



6.

$$y_1' = y_1/(2 + 2t) - 2ty_2, \quad y_2' = y_2/(2 + 2t) + 2ty_1,$$

$a = 0, b = 2$ , точное решение :

$$y_1 = \cos(t^2)\sqrt{1+t}, \quad y_2 = \sin(t^2)\sqrt{1+t}.$$

7.

$$\begin{aligned} y_1' &= -y_2 + t^2 + 6t + 1, \\ y_2' &= y_1 - 3t^2 + 3t + 1, \end{aligned}$$

$a = 0, b = 3$ , точное решение:

$$y_1 = 3t^2 - t - 1 + \cos(t) + \sin(t), \quad y_2 = t^2 + 2 - \cos(t) + \sin(t).$$

8.

$$y_1' = y_2, \quad y_2' = 2y_1^2(1 - 4t^2y_1),$$

$a = 0, b = 5$ , точное решение:

$$y_1 = 1/(1 + t^2), \quad y_2 = -2t/(1 + t^2)^2.$$

9.

$$\begin{aligned} y_1' &= -\sin(t)/(1 + e^{2t})^{1/2} + y_1(y_1^2 + y_2^2 - 1), \\ y_2' &= \cos(t)/(1 + e^{2t})^{1/2} + y_2(y_1^2 + y_2^2 - 1), \end{aligned}$$

$a = 0, b = 5$ , точное решение:

$$y_1 = \cos(t)/(1 + e^{2t})^{1/2}, \quad y_2 = \sin(t)/(1 + e^{2t})^{1/2}.$$

10.

$$\begin{aligned} y_1' &= -\sin(t)/(1 + e^{2t})^{1/2} + y_1(y_1^2 + y_2^2 - 1), \\ y_2' &= \cos(t)/(1 + e^{2t})^{1/2} + y_2(y_1^2 + y_2^2 - 1), \end{aligned}$$

$a = 0, b = 5$ , точное решение:

$$y_1 = \cos(t)/(1 + e^{2t})^{1/2}, \quad y_2 = \sin(t)/(1 + e^{2t})^{1/2}.$$

---

---

## Литература

1. **Бахвалов Н.С., Жидков Н.П., Кобельков Г.М.** Численные методы. — Москва: БИНОМ. Лаб. знаний, 2007.
2. **Самарский А.А., Гулин А.В.** Численные методы. — М.: Наука, 1989.
3. **Справочник** по специальным функциям с формулами, графиками и математическими таблицами/Под ред. М. Абрамовица и И. Стиган. — М.: Наука, 1979.