

В.Б. Андреев

ЧИСЛЕННЫЕ МЕТОДЫ

Часть I

Глава I

Вычислительные методы линейной алгебры

С вычислительной точки зрения в линейной алгебре имеются, если понимать их достаточно широко, две основные задачи:

- 1° решение систем линейных уравнений,
- 2° вычисление собственных значений и собственных векторов матрицы.

Основное внимание в лекциях будет уделено решению первой задачи, да и то при весьма ограничительных предположениях. Вторая задача более трудная — ее мы коснемся менее подробно.

В силу теоремы Кронекера - Капелли система линейных алгебраических уравнений

$$Ax = b$$

разрешима тогда и только тогда, когда ранг матрицы A равен рангу расширенной матрицы $[Ab]$. Это заведомо так, если матрица A квадратная и невырожденная, т.е. $\det A \neq 0$. В этом случае система не только разрешима при любых b , но и имеет единственное решение. (Разрешима однозначно).

Именно этот случай мы и будем изучать.

Методы решения систем линейных алгебраических уравнений делятся на две группы. К первой группе принадлежат так называемые прямые методы — алгоритмы, позволяющие получить решение за конечное число арифметических действий. Сюда относятся известное правило Крамера нахождения решения при помощи определителей, метод исключения Гаусса, метод прогонки — метод решения систем с трехдиагональными матрицами. Существуют и другие методы, из которых отметим метод Холецкого (метод квадратных корней), применяемый к системам с симметричными положительно определенными матрицами, метод вращений и метод отражений.

Вторую группу составляют приближенные методы, в частности, итерационные. В итерационных методах решение системы получается как предел при стремлении числа итераций n к бесконечности. При конечных n , как правило, получаются лишь приближенные решения.

Прямые и итерационные методы имеют свою область применения: если размерность системы не слишком велика, то часто предпочтительнее использовать прямые методы. Итерационные методы выгодны для систем большого порядка. Особенно в случае матриц специального вида.

В настоящем курсе основное внимание будет уделено прямым методам, а итерационных методов коснемся лишь кратко. Более подробно итерационные методы будут изложены во второй части курса "Численные методы", которая читается на четвертом курсе.

§ 1

Метод исключения Гаусса и треугольное (LU) разложение матрицы

1.1 Метод исключения Гаусса

Система $Ax = b$ в развернутой форме имеет вид

$$\sum_{j=1}^n a_{ij}x_j = b_i, \quad i = 1, \dots, n. \quad (1.1)$$

Как известно, метод Гаусса или метод последовательного исключения неизвестных состоит в том, что неизвестные x_j , $j = 1, \dots, n - 1$ последовательно исключаются из соответствующих уравнений системы (1.1), в результате чего она преобразуется к эквивалентной системе с треугольной матрицей

$$\begin{aligned} a_{11}^{(0)}x_1 + a_{12}^{(0)}x_2 + a_{13}^{(0)}x_3 + \dots + a_{1n}^{(0)}x_n &= b_1^{(0)}, \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n &= b_2^{(1)}, \\ \dots & \\ a_{ii}^{(i-1)}x_i + \dots + a_{in}^{(i-1)}x_n &= b_i^{(i-1)}, \\ \dots & \\ a_{nn}^{(n-1)}x_n &= b_n^{(n-1)}, \end{aligned} \quad (1.2)$$

коэффициенты $a_{ij}^{(k)}$ которой и компоненты ее правой части $b_i^{(k)}$ вычисляются по формулам

$$a_{ij}^{(k)} = a_{ij}^{(k-1)} - l_{ik}a_{kj}^{(k-1)}, \quad \begin{aligned} i, j &= k + 1, \dots, n; \\ k &= 1, \dots, n - 1; \end{aligned} \quad (1.3)$$

$$b_i^{(k)} = b_i^{(k-1)} - l_{ik}b_k^{(k-1)}, \quad \begin{array}{l} i, j = k + 1, \dots, n; \\ k = 1, \dots, n - 1; \end{array} \quad (1.4)$$

а

$$l_{ik} = a_{ik}^{(k-1)} / a_{kk}^{(k-1)}, \quad \begin{array}{l} i = k + 1, \dots, n; \\ k = 1, \dots, n - 1; \end{array} \quad (1.5)$$

причем $a_{ij}^{(0)} = a_{ij}$, $b_i^{(0)} = b_i$.

Вычисления по формулам (1.3)-(1.5) называются прямым ходом метода Гаусса. После этого неизвестные x_k последовательно, начиная с x_n , находятся из (1.2) по формулам

$$x_i = \left[b_i^{(i-1)} - \sum_{j=i+1}^n a_{ij}^{(i-1)} x_j \right] / a_{ii}^{(i-1)}, \quad i = n, \dots, 1. \quad (1.6)$$

Вычисления по этим формулам называют обратным ходом метода Гаусса.

Замечание 1.1. В формулах (1.2) при преобразованиях системы (1.1) первое уравнение осталось без изменения. С равным успехом может быть использован и другой вариант исключения, когда первое уравнение (1.1) делится на a_{11} , а вместо (1.2) получается система с единичными коэффициентами при x_j в j -ом уравнении.

Замечание 1.2. Вычисления по формулам (1.5), (1.6), а, следовательно, и по формулам (1.3), (1.4) возможны лишь тогда, когда все числа

$$a_{ii}^{(i-1)} \neq 0, \quad i = 1, \dots, n. \quad (1.7)$$

Необходимые и достаточные условия выполнения (1.7) устанавливаются в доказываемой чуть позже теореме 1.2

1.2 LU разложение матрицы.

Покажем, что метод Гаусса эквивалентен разложению матрицы A системы (1.1) в произведение нижней L и верхней U треугольных матриц с последующим решением вспомогательных систем с этими матрицами. В самом деле, из (1.4) находим, что

$$b_i^{(k-1)} = b_i^{(k-2)} - l_{ik-1}b_{k-1}^{(k-2)}.$$

Подставляя это соотношение в (1.4), получим

$$b_i^{(k)} = b_i^{(k-2)} - l_{ik-1}b_{k-1}^{(k-2)} - l_{ik}b_k^{(k-1)}.$$

Точно так же подставляя сюда выражения для $b_i^{(k-2)}$, а затем для $b_i^{(k-3)}$ и т.д., будем иметь

$$b_i^{(k)} = b_i^{(0)} - l_{i1}b_1^{(0)} - l_{i2}b_2^{(1)} - \dots - l_{ik}b_k^{(k-1)}.$$

Полагая здесь $k = i - 1$ и выражая $b_i^{(0)} = b_i$ через $b_i^{(j)}$, получим

$$b_i = \sum_{j=1}^{i-1} l_{ij} b_j^{(j-1)} + b_i^{(i-1)}. \quad (1.8)$$

Обозначим столбец правой части системы (1.2) через $y = [y_1 \dots y_n]^T$, полагая

$$y_i = b_i^{(i-1)}. \quad (1.9)$$

В этих обозначениях (1.8) переписется так

$$b_i = \sum_{j=1}^{i-1} l_{ij} y_j + y_i, \quad i = 1, \dots, n. \quad (1.10)$$

Обозначим через L нижнюю треугольную матрицу с коэффициентами l_{ij} , вычисляемым по формулам (1.5), и единичной главной диагональю

$$L = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ l_{21} & 1 & 0 & \dots & 0 \\ l_{31} & l_{32} & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & l_{n3} & \dots & 1 \end{bmatrix}. \quad (1.11)$$

Тогда (1.10) можно записать в матричном виде

$$b = Ly. \quad (1.12)$$

Если верхнюю треугольную матрицу системы (1.2) обозначить через U и переписать (1.2) в матричном виде, то будем иметь

$$Ux = y. \quad (1.13)$$

Действуя теперь на левую и правую часть (1.13) невырожденной матрицей L и принимая во внимание (1.12), получим

$$LUx = Ly = b \quad \Rightarrow \quad Ax = b. \quad (1.14)$$

Итак, мы показали, что реализация вычислений по формулам (1.3) и (1.5) прямого хода метода Гаусса эквивалентна разложению матрицы A системы (1.1) в произведение нижней треугольной матрицы с единичной главной диагональю L и верхней треугольной матрицы U . При этом элементы матрицы L вычисляются по формулам (1.5), а элементы матрицы U суть

$$u_{kj} = a_{kj}^{(k-1)} \quad (1.15)$$

и вычисляются по формулам (1.3).

После разложения матрицы A в произведение двух треугольных для отыскания решения системы (1.1) нужно решить две системы с треугольными матрицами — системы (1.12) и (1.13). Решение системы (1.12) заменяет преобразование вектора правой части системы (1.1) по формулам (1.4) прямого хода метода Гаусса. Решение же x системы (1.13) с учетом обозначений (1.9) и (1.15) определяется формулами (1.6) обратного хода метода Гаусса.

Замечание 1.3. Соотношения (1.3) содержат формулы для u_{kj} (1.15) и промежуточные значения, которые тоже нужно запоминать и хранить. Мы сейчас преобразуем эти формулы к такому виду, при котором хранение промежуточных значений не требуется.

Пусть матрица L имеет вид (1.11), т.е. ее элементы

$$l_{ik} = 0 \quad \text{при} \quad k > i, \quad (1.16)$$

а

$$U = \begin{bmatrix} u_{11} & u_{12} & u_{13} & \dots & u_{1n} \\ & u_{22} & u_{23} & \dots & u_{2n} \\ & & u_{33} & \dots & u_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ & & & & u_{nn} \end{bmatrix},$$

т.е.

$$u_{kj} = 0 \quad \text{при} \quad k > j. \quad (1.17)$$

Поскольку $LU = A$, то по правилу умножения матриц находим, что

$$a_{ij} = \sum_{k=1}^n l_{ik} u_{kj}. \quad (1.18)$$

Преобразуем эту формулу двумя способами. В силу (1.11), (1.16)

$$\begin{aligned} \sum_{k=1}^n l_{ik} u_{kj} &= \sum_{k=1}^{i-1} l_{ik} u_{kj} + l_{ii}^{-1} u_{ij} + \sum_{k=i+1}^n \overset{0}{l_{ik}} u_{kj} = \\ &= \sum_{k=1}^{i-1} l_{ik} u_{kj} + u_{ij}, \end{aligned}$$

а в силу (1.17)

$$\begin{aligned} \sum_{k=1}^n l_{ik} u_{kj} &= \sum_{k=1}^{j-1} l_{ik} u_{kj} + l_{ij} u_{jj} + \sum_{k=j+1}^n \overset{0}{l_{ik}} u_{kj} = \\ &= \sum_{k=1}^{j-1} l_{ik} u_{kj} + l_{ij} u_{jj}. \end{aligned}$$

Отсюда и из (1.18) имеем

$$\begin{aligned} u_{ij} &= \left[a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj} \right] \quad i = 1, \dots, n; \quad j = i, \dots, n; \\ l_{ij} &= \frac{1}{u_{jj}} \left[a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right] \quad j = 1, \dots, n; \quad i = j + 1, \dots, n. \end{aligned} \quad (1.19)$$

Очевидно, что реализация формул (1.19) возможна только тогда, когда все $u_{jj} = a_{jj}^{(j-1)}$ в силу (1.15) отличны от нуля (ср. с (1.7)).

Замечание 1.4. Формулы (1.19) устроены так, что нельзя сначала вычислить все u_{ij} , а затем все l_{ij} или наоборот. Можно предложить следующий порядок вычислений по формулам (1.19):

$$\begin{aligned} u_{1j} &= a_{1j}, \quad j = 1, 2, \dots, n; \\ l_{i1} &= a_{i1}/u_{11}, \quad i = 2, 3, \dots, n; \\ u_{2j} &= a_{2j} - l_{21}u_{1j}, \quad j = 2, 3, \dots, n; \\ l_{i2} &= (a_{i2} - l_{i1}u_{12})/u_{22}, \quad i = 3, 4, \dots, n; \end{aligned}$$

и т.д., т.е. чередовать вычисление строк матрицы U и столбцов матрицы L .

После построения матриц L и U решение систем (1.12) и (1.13) с треугольными матрицами находятся по формулам

$$y_i = b_i - \sum_{k=1}^{i-1} l_{ik} y_k, \quad i = 1, 2, \dots, n, \quad (1.20)$$

(вычисления ведутся сверху вниз).

$$x_k = \frac{1}{u_{kk}} \left[y_k - \sum_{j=k+1}^n u_{kj} x_j \right], \quad k = n, n-1, \dots, 1 \quad (1.21)$$

(вычисления ведутся снизу вверх).

Одной из важнейших характеристик любого численного метода является его трудоемкость. Под трудоемкостью метода, предназначенного для решения системы (1.1), обычно понимают число арифметических действий, необходимых для нахождения искомого решения. Часто в трудоемкость метода включают лишь действия умножения и деления, как наиболее трудоемкие операции с точки зрения работы компьютера. Так будем поступать и мы.

Легко видеть, что для вычислений по формулам (1.19) (для получения треугольного разложения) требуется

$$\begin{aligned} Q &= \sum_{i=1}^n \sum_{j=i}^n (i-1) + \sum_{j=1}^n \sum_{i=j+1}^n j = \sum_{i=1}^n [(i-1)(n-i+1) + i(n-i)] = \\ &= 2 \sum_{i=1}^n [(n+1)i - i^2] - n(n+1) = n(n+1)^2 - \frac{n(n+1)(2n+1)}{3} - n(n+1) = (1.22) \\ &= \frac{n(n^2-1)}{3} = \frac{n^3}{3} + O(n) \approx \frac{n^3}{3} \end{aligned}$$

действий умножения и деления.

Для вычислений по формулам (1.20) и (1.21) имеем соответственно

$$\overset{\circ}{q} = \sum_{i=1}^n (i-1) = \frac{n(n-1)}{2} \quad \text{и} \quad \bar{q} = \sum_{k=1}^n (n-k+1) = \frac{n(n+1)}{2},$$

т.е. общее число действий для решения систем (1.12) и (1.13) по формулам (1.20), (1.21) есть

$$q = \overset{\circ}{q} + \bar{q} = n^2. \quad (1.23)$$

Замечание 1.5. Из формул (1.22) и (1.23) следует, что при больших n основной объем работы, которую нужно выполнить для решения системы (1.1) описанным методом, падает на преобразование коэффициентов матрицы системы, т.е. на построение треугольного разложения, в то время как преобразование вектора правой части (решение системы (1.12)) и на отыскание самого решения трудозатраты сравнительно невелики. В связи с этим при больших n решение нескольких систем с различными правыми частями и одной и той же матрицей оказывается по трудоемкости практически таким же как и решение одной системы.

Выясним теперь условия, при которых вычисления по формулам (1.19)-(1.21) возможны, т.е. все u_{jj} отличны от нуля.

Теорема 1.1. Пусть A — невырожденная матрица, L — нижняя треугольная матрица с единичной главной диагональю, а U — невырожденная верхняя треугольная матрица. Тогда, если $A = LU$, то это представление единственно.

Для доказательства теоремы 1.1 нам потребуется

Лемма 1.1. Произведение нижних (верхних) треугольных матриц есть нижняя (верхняя) треугольная матрица. Обратная к невырожденной нижней (верхней) треугольной матрице есть нижняя (верхняя) треугольная матрица.

Упражнение 1.1. Доказать лемму 1.1.

Доказательство теоремы 1.1. Пусть $A = L_1U_1 = L_2U_2$. Тогда

$$L_2 = L_1U_1U_2^{-1} \quad \text{и} \quad L_1^{-1}L_2 = U_1U_2^{-1}.$$

Слева стоит произведение нижних треугольных матриц, а справа — верхних. Поэтому произведение есть диагональная матрица D , т.е. $L_1^{-1}L_2 = D$. Отсюда находим, что $L_2 = L_1D$. Поскольку главные диагонали L_1 и L_2 единичные, то главная диагональ L_1D совпадает с главной диагональю D , следовательно, $D = I$. Отсюда $L_1 = L_2$ и $U_1 = U_2$. Теорема доказана.

Теорема 1.2. Пусть A — квадратная невырожденная матрица, L — нижняя треугольная матрица с единичной главной диагональю, а U — невырожденная верхняя треугольная матрица. Разложение $A = LU$ существует тогда и только тогда, когда все угловые миноры матрицы A отличны от нуля.

Напомним, что угловыми минорами матрицы A называются величины

$$\Delta_1 = a_{11}, \quad \Delta_2 = \det \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \quad \dots, \quad \Delta_n = \det[A].$$

Доказательство. 1°. (Необходимость)

Пусть разложение $A = LU$ существует. Тогда по теореме 1.1 оно единственно. Представим матрицы A , L и U в блочном виде

$$A = \begin{bmatrix} A_m & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad L = \begin{bmatrix} L_m & 0 \\ L_{21} & L_{22} \end{bmatrix}, \quad U = \begin{bmatrix} U_m & U_{12} \\ 0 & U_{22} \end{bmatrix},$$

где A_m , L_m , U_m и A_{22} , L_{22} , U_{22} — квадратные матрицы размерностей $m \times m$ и $(n - m) \times (n - m)$ соответственно, а $m < n$ — произвольное число. Разложение $A = LU$ в блочном представлении имеет вид

$$\begin{bmatrix} A_m & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} L_m & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} U_m & U_{12} \\ 0 & U_{22} \end{bmatrix} = \begin{bmatrix} L_m U_m & L_m U_{12} \\ L_{21} U_m & L_{21} U_{12} + L_{22} U_{22} \end{bmatrix} \quad (1.24)$$

Отсюда следует, что

$$A_m = L_m U_m. \quad (1.25)$$

Поскольку матрица U треугольная и невырожденная, то все ее диагональные элементы отличны от нуля. Поэтому невырождена и треугольная матрица U_m . Тем самым

$$\Delta_m = \det[A_m] = \det[L_m] \det[U_m] = u_{11} \dots u_{mm} \neq 0$$

при $m = 1, \dots, n$.

2°. (Достаточность) Пусть теперь $\Delta_1 \Delta_2 \dots \Delta_n \neq 0$. Для доказательства существования треугольного разложения воспользуемся методом полной математической индукции по порядку системы n . При $n = 1$ матрица $A = a_{11} = \Delta_1 \neq 0$, матрица $L = 1$

и поэтому $U = u_{11} = a_{11} = \det U \neq 0$. Существование искомого разложения при $n = 1$ доказано.

Пусть A_k — матрица порядка k и разложение $A_k = L_k U_k$ существует с $\det U_k \neq 0$ при $k = 1, \dots, m-1$. Докажем, что существует и $A_m = L_m U_m$, причем $\det U_m \neq 0$. Пусть $a_{\cdot m} = [a_{1m} \dots a_{m-1m}]^T$ — столбец, $a_{m\cdot} = [a_{m1} \dots a_{mm-1}]$ — строка и разложение A_m будем искать в виде

$$\begin{aligned} A_m &= \begin{bmatrix} A_{m-1} & a_{\cdot m} \\ a_{m\cdot} & a_{mm} \end{bmatrix} = \begin{bmatrix} L_{m-1} & 0 \\ l_{m\cdot} & 1 \end{bmatrix} \begin{bmatrix} U_{m-1} & u_{\cdot m} \\ 0 & u_{mm} \end{bmatrix} = \\ &= \begin{bmatrix} L_{m-1} U_{m-1} & L_{m-1} u_{\cdot m} \\ l_{m\cdot} U_{m-1} & l_{m\cdot} u_{\cdot m} + u_{mm} \end{bmatrix}. \end{aligned}$$

Отсюда следует, что неизвестный столбец $u_{\cdot m}$, неизвестная строка $l_{m\cdot}$ и u_{mm} определяются следующими соотношениями:

$$\begin{aligned} L_{m-1} u_{\cdot m} &= a_{\cdot m} \\ l_{m\cdot} U_{m-1} &= a_{m\cdot} \quad \Rightarrow \quad U_{m-1}^T l_{m\cdot}^T = a_{m\cdot}^T, \\ u_{mm} &= a_{mm} - l_{m\cdot} u_{\cdot m}. \end{aligned} \tag{1.26}$$

Поскольку L_{m-1} и U_{m-1} невырожденные, из первого и второго соотношений (1.26) можно найти $u_{\cdot m}$ и $l_{m\cdot}$, соответственно, после чего третье соотношение дает u_{mm} . Существование разложения (1.24) для m доказано. Осталось доказать, что $\det U_m \neq 0$. Но с учетом (1.25)

$$0 \neq \Delta_m = \det A_m = \det U_m,$$

что и требовалось доказать. Теорема полностью доказана.

§ 2

Ленточные методы

2.1 Метод прогонки

Рассмотрим систему линейных алгебраических уравнений

$$Ax = b, \tag{2.1}$$

матрица которой является трехдиагональной. Запишем эту систему в развернутом виде. Пусть

$$\begin{aligned} b_1x_1 + c_1x_2 &= d_1, \\ a_2x_1 + b_2x_2 + c_2x_3 &= d_2, \\ \dots & \\ a_ix_{i-1} + b_ix_i + c_ix_{i+1} &= d_i, \\ \dots & \\ a_nx_{n-1} + b_nx_n &= d_n. \end{aligned} \tag{2.2}$$

Алгоритм метода прогонки — метода решения системы (2.2) — состоит в следующем (см. курс "Введение в численные методы", но, может быть, с другими обозначениями!)

а) Нахождение прогоночных коэффициентов (прямая прогонка) по формулам

$$\begin{aligned} \alpha_i &= -c_i/\gamma_i, \quad i = 1, 2, \dots, n-1, \\ \gamma_i &= b_i + a_i\alpha_{i-1}, \quad i = 2, \dots, n, \quad \gamma_1 = b_1, \\ \beta_i &= (d_i - a_i\beta_{i-1})/\gamma_i, \quad i = 2, \dots, n, \quad \beta_1 = d_1/b_1. \end{aligned} \tag{2.3}$$

б) Нахождение самого решения (обратная прогонка)

$$x_i = \alpha_ix_{i+1} + \beta_i, \quad i = n-1, \dots, 1, \quad x_n = \beta_n. \tag{2.4}$$

Из (2.3),(2.4) следует, что общее число умножений и делений при вычислении коэффициентов α_i и γ_i

$$Q = 2(n - 1), \quad (2.5)$$

а при вычислении коэффициентов β_i и решения x_i

$$q = 3(n - 1). \quad (2.6)$$

Сравнение (2.5),(2.6) с (1.22), (1.23) свидетельствуют о том, что прогонка существенно менее трудоемка по сравнению с общим методом Гаусса. Связано это с тем, что мы явным образом воспользовались тем, что значительная часть элементов матрицы A равна нулю.

2.2 Ленточные матрицы

Определение 2.1. Матрица A называется ленточной с полушириной ленты p , если ее элементы $a_{ij} = 0$ при $|i - j| > p$, но существует по крайней мере один элемент $a_{ij} \neq 0$ при $|i - j| = p$

Пример 2.1. Для диагональной матрицы $a_{ij} = 0$ при $|i - j| > 0$ и, следовательно, ее полуширина равна нулю. Ее лента состоит из одной диагонали и ширина равна 1. Условно диагональную матрицу можно изобразить как на рис. 1.

$$\begin{bmatrix} * & & & & & \\ & * & & & & \\ & & * & & & \\ & & & * & & \\ & & & & * & \\ & & & & & * \end{bmatrix}$$

Рис. 1.

Пример 2.2. Полная матрица имеет $2n - 1$ диагоналей. Это и есть ширина ее ленты, а полуширина будет $p = n - 1$.

$$\begin{bmatrix} * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \end{bmatrix}$$

Рис. 2.

Пример 2.3. На рис. 3 изображена матрица с полушириной $p = 1$.

$$\begin{bmatrix} * & * & & & & & \\ * & * & * & & & & \\ & * & * & * & & & \\ & & * & * & * & & \\ & & & * & * & * & \\ & & & & * & * & \\ & & & & & * & * \end{bmatrix}$$

Рис. 3.

Матрицы такой структуры называются трехдиагональными. Ширина ленты 3.

Пример 2.4. Матрица, изображенная на рис.4, имеет полуширину $p = 1$ и ширину ленты 2. Матрицы такой структуры называются трапециевидными. Изображенная матрица также называется правой ленточной.

$$\begin{bmatrix} * & * & & & & & \\ & * & * & & & & \\ & & * & * & & & \\ & & & * & * & & \\ & & & & * & * & \\ & & & & & * & * \\ & & & & & & * \end{bmatrix}$$

Рис. 4.

Пример 2.5. У матрицы на рис. 5

$$\begin{bmatrix} * & 0 & * & & & & \\ 0 & * & 0 & * & & & \\ * & 0 & * & 0 & * & & \\ & * & 0 & * & 0 & * & \\ & & * & 0 & * & 0 & \\ & & & * & 0 & * & \end{bmatrix}$$

Рис. 5.

полуширина $p = 2$, ширина 5 и всего три ненулевых диагонали.

Определение 2.2. Рисунок, на котором (звездочками) отмечены позиции, где только и могут располагаться ненулевые элементы матрицы A , называется ее портретом.

2.3 Ленточный вариант треугольного разложения

Модифицируем алгоритм исключения Гаусса на случай ленточных матриц, т.е. заранее отбросим те вычисления, которые заведомо приводят к нулевым элементам. Это

позволит нам сэкономить в трудозатратах на решение системы. Обратимся сразу к варианту, основанному на треугольном разложении матрицы A .

Нам потребуется

Лемма 2.1. *Если полуширина матрицы A равна p , то в треугольном разложении $A = LU$ полуширина L (U) не больше p .*

Доказательство. Пусть

$$a_{ij} = 0 \quad \text{при} \quad |i - j| > p. \quad (2.7)$$

Докажем, что

$$l_{ij} = 0 \quad \text{при} \quad i - j > p. \quad (2.8)$$

Для доказательства применим метод полной математической индукции по номерам столбцов матрицы L . При $j = 1$ из (1.19) находим, что

$$l_{i1} = a_{i1}/a_{11}, \quad i = 2, \dots, n.$$

Отсюда с учетом (2.7) приходим к (2.8) с $j = 1$, т.е. $l_{i1} = 0$ при $i - 1 > p$. Пусть теперь утверждение (2.8) верно для столбцов матрицы L с номерами $k = 1, 2, \dots, j - 1$, т.е.

$$l_{ik} = 0 \quad \text{при} \quad i - k > p, \quad k = 1, 2, \dots, j - 1. \quad (2.9)$$

Докажем справедливость (2.8) для j -ого столбца. Пусть

$$i - j > p. \quad (2.10)$$

Тогда в силу (1.19) и (2.7)

$$l_{ij} = \frac{1}{u_{jj}} \left[a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right] = -\frac{1}{u_{jj}} \sum_{k=1}^{j-1} l_{ik} u_{kj}. \quad (2.11)$$

Оценим разность индексов $i - k$ у первых сомножителей под знаком суммы в (2.11). С учетом (2.10) и (2.11) будем иметь

$$i - k > p + j - k \geq p + 1.$$

Но тогда в силу (2.9) первые сомножители в сумме (2.11) обращаются в нуль и соотношение (2.8) установлено. Лемма доказана.

Преобразуем формулы (1.19)-(1.21) на случай ленточной матрицы A . Сначала выясним, для каких значений индексов i и j нужно проводить вычисления по формулам (1.19). Так как в силу леммы 2.1

$$u_{ij} = 0 \quad \text{при} \quad j - i > p, \quad (2.12)$$

то ненулевые элементы u_{ij} могут быть лишь при $j-i \leq p$, т.е. при $j \leq p+i$. Аналогично

$$l_{ij} = 0 \quad \text{при} \quad i - j > p \quad (2.13)$$

и, следовательно, ненулевые элементы l_{ij} могут быть только при $i - j \leq p$, т.е. при $i \leq p + j$. Отсюда

$$\begin{aligned} u_{ij} \neq 0 & \quad \begin{aligned} i &= 1, \dots, n, \\ j &= i, \dots, \min[n, p+i], \end{aligned} \\ l_{ij} \neq 0 & \quad \begin{aligned} j &= 1, \dots, n, \\ i &= j+1, \dots, \min[n, p+j]. \end{aligned} \end{aligned}$$

Теперь преобразуем суммы в (1.19). В силу (2.12) ненулевые слагаемые в суммах (1.19) могут быть только при $j - k \leq p$, т.е. при

$$k \geq j - p.$$

а в силу (2.13) — только при $i - k \leq p$, т.е. при

$$k \geq i - p.$$

Объединяя эти неравенства и принимая во внимание, что k — натуральное, будем иметь

$$k \geq \max[1, i - p, j - p].$$

Поскольку в формулах для u_{ij} индексы i, j подчинены ограничению $j \geq i$, то в этих формулах

$$k \geq \max[1, j - p].$$

В формулах же для l_{ij} наоборот $i > j$ и поэтому в них

$$k \geq \max[1, i - p].$$

С учетом сказанного, для ленточной матрицы A с полушириной p формулы (1.19) принимают вид

$$\begin{aligned} u_{ij} &= \left[a_{ij} - \sum_{k=\max[1, j-p]}^{i-1} l_{ik} u_{kj} \right] & \begin{aligned} i &= 1, \dots, n, \\ j &= i, \dots, \min[n, p+i], \end{aligned} \\ l_{ij} &= \frac{1}{u_{jj}} \left[a_{ij} - \sum_{k=\max[1, i-p]}^{j-1} l_{ik} u_{kj} \right] & \begin{aligned} j &= 1, \dots, n, \\ i &= j+1, \dots, \min[n, p+j]. \end{aligned} \end{aligned} \quad (2.14)$$

Преобразуем формулы (1.20) и (1.21). В формулах (1.20) в силу (2.13) $i - k \leq p$, а в формулах (1.21) в силу (2.12) $j - k \leq p$. Поэтому

$$\begin{aligned} y_i &= b_i - \sum_{k=\max[1, i-p]}^{i-1} l_{ik} y_k, & i &= 1, \dots, n, \\ x_k &= \frac{1}{u_{kk}} \left[y_k - \sum_{j=k+1}^{\min[n, k+p]} u_{kj} x_j \right], & k &= n, \dots, 1. \end{aligned} \quad (2.15)$$

2.4 Оценка трудоемкости

Оценим трудоемкость LU -разложения ленточной матрицы A , имеющей полуширину p , и трудоемкость решения системы (2.1) с такой матрицей. Для этого подсчитаем число умножений и делений, необходимых для реализации формул (2.14) и (2.15). На рис. 6 изображены портреты матриц L и U

$$\begin{array}{c}
 1 \\
 p \\
 n-p \left\{ \begin{array}{l} p+1 \\ n \end{array} \right. \\
 \\
 \\
 L
 \end{array}
 \begin{bmatrix}
 |1 & & & & \\
 * & 1 & & & \\
 * & * & \underline{1} & & \\
 * & * & * & 1 & \\
 & * & * & * & 1 \\
 & & * & * & * & 1
 \end{bmatrix}
 \qquad
 \begin{array}{c}
 1 \\
 \\
 \\
 \\
 \\
 \\
 \\
 \\
 \\
 U
 \end{array}
 \begin{bmatrix}
 * & * & * & * & & \\
 & * & * & * & * & \\
 & & * & * & * & \\
 & & & * & * & * \\
 & & & & * & * \\
 & & & & & *
 \end{bmatrix}
 \begin{array}{c}
 \\
 \\
 \\
 \\
 \overbrace{p \ p+1 \ n}^{n-p} \\
 \\
 \\
 \\
 \\
 \\
 \end{array}$$

Рис. 6

Перепишем формулы (2.14), отделив формулы для элементов u_{ij} и l_{ij} , обведенных треугольниками на рис. 6.

$$u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj}, \quad i = 1, \dots, p, \quad j = i, \dots, p, \quad (2.16)$$

$$u_{ij} = a_{ij} - \sum_{k=j-p}^{i-1} l_{ik} u_{kj}, \quad j = p+1, \dots, n, \quad i = j-p, \dots, j, \quad (2.17)$$

$$l_{ij} = \frac{1}{u_{jj}} \left[a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right], \quad j = 1, \dots, p, \quad (2.18)$$

$$i = j+1, \dots, p,$$

$$l_{ij} = \frac{1}{u_{jj}} \left[a_{ij} - \sum_{k=i-p}^{j-1} l_{ik} u_{kj} \right], \quad i = p+1, \dots, n, \quad (2.19)$$

$$j = i-p, \dots, i-1.$$

Подсчитаем число умножений и делений, необходимых для реализации формул (2.16)-(2.19). Из сравнения (2.16), (2.18) с (1.19) следует, что эти формулы при $p = n$ совпадают. Поэтому, с учетом (1.22)

$$Q_{16}(U_p) + Q_{18}(L_p) = \frac{p(p^2 + 1)}{3}. \quad (2.20)$$

Здесь Q — трудоемкость, нижний индекс — номер формулы, а аргумент — объект вычислений.

Далее

$$\begin{aligned} Q_{17}(u_{ij}) &= i - 1 - j + p + 1 = i - j + p, \\ Q_{17}(U) &= \sum_{j=p+1}^n \sum_{i=j-p}^j (i - j + p), \\ Q_{19}(l_{ij}) &= j - 1 - i + p + 1 + 1 = j - i + p + 1, \\ Q_{19}(L) &= \sum_{i=p+1}^n \sum_{j=i-p}^{i-1} (j - i + p + 1) = \sum_{j=p+1}^n \sum_{i=j-p}^{j-1} (i - j + p + 1), \end{aligned}$$

и, следовательно,

$$\begin{aligned} Q_{17}(U) + Q_{19}(L) &= \sum_{j=p+1}^n \left[p + \sum_{i=j-p}^{j-1} (i - j + p + i - j + p + 1) \right] = \\ &= 2 \sum_{j=p+1}^n \left[p + \sum_{k=1}^{p-1} k \right] = (n - p)p(p - 1). \end{aligned}$$

Отсюда и из (2.20) находим, что общее число умножений и делений при построении LU -разложения есть

$$Q = \frac{p(p+1)}{3}(3n - 2p - 1) = n(p^2 + p) - \frac{2}{3}p^3 + O(p^2). \quad (2.21)$$

Замечание 2.1. Полуширина полной матрицы $p = n - 1$. Подставляя это значение p в найденное выражение, получим выражение для трудоемкости треугольного разложения, совпадающее с (1.22). Полагая же здесь $p = 1$, получим (2.5).

Обратимся к формулам (2.15).

$$\begin{aligned} Q_{15}(y_i) &= \begin{bmatrix} i - 1, & i = 1, \dots, p, \\ i - 1 - i + p + 1 = p, & i = p + 1, \dots, n, \end{bmatrix} \\ Q_{15}(y) &= \sum_{i=1}^n Q_{15}(y_i) = \frac{p(p-1)}{2} + p(n-p) = \frac{p(2n-p-1)}{2}, \\ Q_{15}(x) &= \frac{p(2n-p-1)}{2} + n \end{aligned}$$

и следовательно

$$q = Q_{15}(x) + Q_{15}(y) = p(2n - p - 1) + n = (2p + 1)n - p(p + 1) \quad (2.22)$$

(ср. с (1.23) при $p = n - 1$ и (2.6) при $p = 1$).

Проанализируем формулы (2.21), (2.22). Рассмотрим три случая.

1°. $p = O(n)$, например, $p = \alpha n$, $\alpha < 1$. Тогда

$$Q \approx \alpha^2 n^3 - \frac{2}{3} \alpha^3 n^3 = \alpha^2 \left(1 - \frac{2\alpha}{3}\right) n^3.$$

Легко проверить, что при $0 < \alpha < 1$ коэффициент при n^3 меньше $1/3$.

2°. $p = o(n)$, но $p \rightarrow \infty$ при $n \rightarrow \infty$. В этом случае

$$Q \approx p^2 n, \quad q \approx 2pn.$$

В частности, при $p = \sqrt{n}$

$$Q \approx n^2, \quad q = 2n^{3/2}.$$

3°. $p = \text{const}$ ($p = 1, 2, \dots$).

$$Q \approx p(p+1)n, \quad q \approx (2p+1)n.$$

При $p = 1$ $Q < q$, при $p \geq 2$ $Q > q$.

2.5 LU - разложение для трехдиагональной матрицы.

Линейная система (2.1) с трехдиагональной матрицей A представляет особый интерес в силу того, что часто встречается в приложениях. Получим формулы LU - разложения для этого случая, не апеллируя к (2.14), (2.15). Запишем систему (2.1) в виде (2.2). Найдём такое LU - разложение матрицы A , в котором U (а не L !) имеет единичную главную диагональ (см. замечание 1.1). Матрица системы (2.2) имеет вид

$$A = \begin{bmatrix} b_1 & c_1 & 0 & \dots & 0 \\ a_2 & b_2 & c_2 & \dots & 0 \\ 0 & a_3 & b_3 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & b_n \end{bmatrix} \quad (2.23)$$

Пусть

$$L = \begin{bmatrix} \gamma_1 & 0 & 0 & \dots & 0 \\ \delta_2 & \gamma_2 & 0 & \dots & 0 \\ 0 & \delta_3 & \gamma_3 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \gamma_n \end{bmatrix}, \quad U = \begin{bmatrix} 1 & -\alpha_1 & 0 & \dots & 0 \\ 0 & 1 & -\alpha_2 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}.$$

Тогда

$$LU = \begin{bmatrix} \gamma_1 & -\gamma_1 \alpha_1 & & & \\ \delta_2 & -\delta_2 \alpha_1 + \gamma_2 & -\gamma_2 \alpha_2 & & \\ & \delta_3 & -\delta_3 \alpha_2 + \gamma_3 & -\gamma_3 \alpha_3 & \\ & & \delta_4 & -\delta_4 \alpha_3 + \gamma_4 & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

Сравнивая эту формулу с матрицей (2.23), находим формулы для элементов матриц U и L :

$$\begin{aligned} \alpha_i &= -c_i/\gamma_i, \quad i = 1, \dots, n-1, \\ \gamma_i &= b_i + \delta_i\alpha_{i-1}, \quad i = 2, \dots, n, \quad \gamma_1 = b_1, \\ \delta_i &= a_i, \quad i = 2, \dots, n. \end{aligned} \tag{2.24}$$

Обратимся к решению системы (2.2) с использованием LU - разложения. Пусть $Ux = \beta$. Тогда $L\beta = d$. Записывая последнюю систему в развернутом виде, будем иметь

$$\begin{aligned} \gamma_1\beta_1 &= d_1, \\ \delta_2\beta_1 + \gamma_2\beta_2 &= d_2, \\ \dots\dots\dots & \\ \delta_i\beta_{i-1} + \gamma_i\beta_i &= d_i, \\ \dots\dots\dots & \\ \delta_n\beta_{n-1} + \gamma_n\beta_n &= d_n. \end{aligned}$$

Отсюда

$$\beta_i = (d_i - \delta_i\beta_{i-1})/\gamma_i, \quad i = 2, \dots, n, \quad \beta_1 = \frac{d_1}{\gamma_1}. \tag{2.25}$$

Аналогично находим, что решение системы $Ux = \beta$ определяется формулами

$$x_i = \alpha_i x_{i+1} + \beta_i, \quad i = n-1, \dots, 1, \quad x_n = \beta_n. \tag{2.26}$$

Сравнение формул (2.24) - (2.26) с формулами (2.3),(2.4) показывает, что метод прогонки является не чем иным как частным случаем ленточного варианта метода исключения Гаусса.

Определение 2.3. Говорят, что матрица A ленточная с лентой нижней полуширины p_1 и верхней полуширины p_2 , если $a_{ij} = 0$ при $i - j > p_1$ и $j - i > p_2$.

Упражнение 2.1. Построить модификацию формул (2.14), (2.15) для случая матрицы с разной полушириной.

§ 3

Методы Холецкого и блочного исключения. Вычисление обратной матрицы

3.1 Метод Холецкого (квадратных корней)

Вновь обратимся к системе

$$Ax = b. \quad (3.1)$$

На этот раз будем предполагать, что матрица A симметрична и положительно определена, т.е.

$$A = A^T \quad \text{и} \quad A > 0. \quad (3.2)$$

Последнее означает, что квадратичная форма $x^T Ax > 0$ для любого ненулевого вектора x . Напомним, что симметричная матрица имеет только действительные собственные значения, а положительно определенная — только положительные. В силу критерия Сильвестра необходимым и достаточным условием положительной определенности матрицы A является положительность всех ее угловых миноров $\Delta_i > 0$, $i = \dots, n$.

Построим алгоритм решения системы (3.1), который использует свойства (3.2) матрицы A . Это будет метод Холецкого. Основой метода Холецкого является

Теорема 3.1. *Если $A = A^T > 0$, то существует единственное разложение*

$$A = LL^T, \quad (3.3)$$

где L — нижняя треугольная матрица с положительными диагональными элементами.

Определение 3.1. Разложение (3.3) называется разложением Холецкого, а матрица L — множителем Холецкого.

Доказательство теоремы. Сначала докажем единственность. Пусть существуют два разложения

$$A = L_1 L_1^T = L_2 L_2^T.$$

Обращая матрицы L_2 и L_1^T , будем иметь

$$L_2^{-1} L_1 = L_2^T (L_1^T)^{-1}.$$

Принимая во внимание, что $(AB)^T = B^T A^T$, а для невырожденных матриц $(AB)^{-1} = B^{-1} A^{-1}$, найдем, что

$$(L_1^{-1} L_2)^{-1} = L_2^{-1} L_1 = L_2^T (L_1^T)^{-1} = (L_1^{-1} L_2)^T. \quad (3.4)$$

В силу леммы 1.1 обратная к нижней треугольной матрице есть нижняя треугольная матрица и произведение таких матриц есть снова нижняя треугольная матрица. Из сказанного следует, что в левой части (3.4) стоит нижняя треугольная матрица, а справа — верхняя. Равенство (3.4) возможно только тогда, когда обе матрицы диагональные. Но диагональная матрица совпадает со своей транспонированной. Поэтому из (3.4) следует, что

$$(L_1^{-1} L_2)^{-1} = L_1^{-1} L_2 = D.$$

Это соотношение утверждает, что диагональная матрица D совпадает со своей обратной, что возможно только в том случае, если у этой матрицы диагональными элементами являются числа ± 1 . Поскольку $L_2 = L_1 D$, а диагональные элементы L_1 и L_2 положительны, то диагональные элементы D тоже должны быть положительны, т.е. $D \equiv I$ и, следовательно, $L_2 = L_1$. Единственность доказана.

Построим теперь формулы для вычисления элементов L , откуда и будет следовать существование. Так как $a_{ij} = a_{ji}$, а $l_{ij} = 0$ при $i < j$, то будем считать, что

$$i \geq j.$$

Тогда

$$\begin{aligned} a_{ij} &= \sum_{k=1}^n l_{ik} l_{kj}^T = \sum_{k=1}^n l_{ik} l_{jk} = \sum_{k=1}^{j-1} l_{ik} l_{jk} + l_{ij} l_{jj} + \sum_{k=j+1}^n l_{ik} l_{jk} = \\ &= \sum_{k=1}^{j-1} l_{ik} l_{jk} + l_{ij} l_{jj}. \end{aligned}$$

При $i = j$ находим, что

$$l_{jj} = \sqrt{a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2}, \quad j = 1, \dots, n. \quad (3.5)$$

Далее,

$$l_{ij} = \frac{1}{l_{jj}} \left[a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} \right], \quad i = j+1, \dots, n, \quad (3.6)$$

$$j = 1, \dots, n-1$$

Вычисления можно вести по столбцам $j = 1, \dots, n$ для $i = j+1, \dots, n$.

$$j = 1: \quad l_{11} = \sqrt{a_{11}}, \quad l_{i1} = \frac{a_{i1}}{l_{11}}, \quad i = 2, \dots, n$$

$$j = 2: \quad l_{22} = \sqrt{a_{22} - l_{21}^2}, \quad l_{i2} = \frac{1}{l_{22}} [a_{i2} - l_{i1} l_{21}], \quad i = 3, \dots, n$$

и т.д.

Осталось доказать, что все l_{jj} положительны, т.е. положительны подкоренные выражения. Докажем, что

$$l_{jj} = \sqrt{\Delta_j / \Delta_{j-1}}, \quad \text{где } \Delta_0 = 1.$$

Пусть, как раньше,

$$A = \begin{bmatrix} A_j & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad L = \begin{bmatrix} L_j & 0 \\ L_{21} & L_{22} \end{bmatrix}.$$

Тогда

$$A_j = L_j L_j^T.$$

Отсюда

$$\Delta_j = \det A_j = (\det L_j)^2 = \left(\prod_{k=1}^j l_{kk} \right)^2.$$

Аналогично

$$\Delta_{j-1} = \left(\prod_{k=1}^{j-1} l_{kk} \right)^2$$

и, следовательно,

$$l_{jj}^2 = \Delta_j / \Delta_{j-1} > 0, \quad j = 1, \dots, n.$$

Теорема доказана.

Упражнение 3.1. Показать, что для реализации формул (3.5), (3.6) при всех i и j требуется

$$Q = \frac{n(n+1)(n+2)}{6} \approx \frac{n^3}{6} \quad (3.7)$$

операций умножения, деления и извлечения корня.

Замечание 3.1. Из (3.7) следует, что разложение Холецкого в два раза более экономично, чем треугольное разложение.

Обратимся теперь к решению системы (3.1). Поскольку $Ax = LL^T x = b$, то, полагая $L^T x = y$, получим $Ly = b$. При этом

$$\begin{aligned} y_i &= \frac{1}{l_{ii}} \left[b_i - \sum_{k=1}^{i-1} l_{ik} y_k \right], \quad i = 1, \dots, n, \\ x_i &= \frac{1}{l_{ii}} \left[y_i - \sum_{k=i+1}^n l_{ki} x_k \right], \quad i = n, \dots, 1. \end{aligned} \quad (3.8)$$

Замечание 3.2. Очень часто используется модификация разложения Холецкого, называемая LDL^T -разложением. Суть ее в том, что вместо разложения (3.3) строится разложение матрицы A вида

$$A = LDL^T,$$

где L — попрежнему нижняя треугольная матрица, но в отличие от (3.3) ее диагональные элементы равны 1, а D — диагональная матрица. Достоинство LDL^T -разложения состоит в том, что при его вычислении не требуется находить квадратные корни, а потому оно существует не только для положительно определенных матриц. Условием существования такого разложения для симметричной матрицы является отличие от нуля всех ее угловых миноров.

Упражнение 3.2. Построить формулы для вычисления элементов матриц L и D в LDL^T -разложении.

Ответ.

$$\begin{aligned} d_j &= a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2 d_k, \quad j = 1, 2, \dots, n. \\ l_{ij} &= \left[a_{ij} - \sum_{k=1}^{j-1} l_{ik} d_k l_{jk} \right] / d_j, \quad i = j + 1, \dots, n, \quad j = 1, \dots, n - 1. \end{aligned}$$

3.2 Ленточный вариант метода Холецкого

Как и в случае треугольного разложения можно построить разложение Холецкого в ленточном варианте. Пусть A — симметричная положительно определенная ленточная матрица с полушириной p . Справедлива

Лемма 3.1. Если полуширина матрицы $A = A^T > 0$ равна p , то и полуширина множителя Холецкого не больше p .

На доказательстве этой леммы мы не останавливаемся, ибо оно почти дословно повторяет доказательство аналогичной леммы 2.1.

В силу леммы 3.1 ненулевыми элементами l_{ik} матрицы L могут быть только те, у которых индексы подчинены условиям

$$j - p \leq k \leq j \quad (l_{jk} \neq 0). \quad (3.9)$$

Поэтому формула (3.5) принимает вид

$$l_{jj} = \sqrt{a_{jj} - \sum_{k=\max[1, j-p]}^{j-1} l_{jk}^2}, \quad j = 1, \dots, n. \quad (3.10)$$

В силу (3.9) в формулах (3.6) $i - p \leq k \leq i$ и $i - p \leq j \leq i$. Поэтому формулы (3.6) принимают вид

$$l_{ij} = \frac{1}{l_{jj}} \left[a_{ij} - \sum_{k=\max[1, i-p]}^{j-1} l_{ik} l_{jk} \right], \quad \begin{array}{l} i = j + 1, \dots, \min[n, p + j], \\ j = 1, \dots, n. \end{array} \quad (3.11)$$

Ленточный вариант разложения Холецкого построен.

Для отыскания решения системы (3.1) нужно еще преобразовать формулы (3.8). Они принимают вид

$$\begin{aligned} y_i &= \frac{1}{l_{ii}} \left[b_i - \sum_{k=\max[1, i-p]}^{i-1} l_{ik} y_k \right], \quad i = 1, \dots, n, \\ x_i &= \frac{1}{l_{ii}} \left[y_i - \sum_{k=i+1}^{\min[n, i+p]} l_{ki} x_k \right], \quad i = n, \dots, 1. \end{aligned} \quad (3.12)$$

Эти формулы полностью аналогичны формулам (??).

Упражнение 3.3. Подсчитать число действий умножения, деления и извлечения корня, необходимых для реализации формул (3.10)-(3.11).

Ответ.

$$Q = \frac{(p+1)(p+2)(3n-2p)}{6}.$$

3.3 Метод блочного исключения (метод частичного исключения неизвестных)

В методе исключения Гаусса из системы (3.1) последовательно исключаются неизвестные — компоненты вектора $x^T = [x_1, x_2, \dots, x_n]$. В ряде случаев бывает полезным процедуру исключения неизвестных произвести блочно. Пусть

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad (3.13)$$

где A_{11} — квадратная невырожденная матрица размеров $m \times m$, а b_1 и x_1 — m -мерные векторы. С учетом (3.13) система (3.1) принимает вид

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

или после блочного перемножения

$$\begin{aligned} A_{11}x_1 + A_{12}x_2 &= b_1, \\ A_{21}x_1 + A_{22}x_2 &= b_2. \end{aligned} \tag{3.14}$$

Из первого уравнения (3.14) находим, что

$$x_1 = A_{11}^{-1}(b_1 - A_{12}x_2). \tag{3.15}$$

Подставляя это представление x_1 во второе уравнение (3.14), получим

$$A_{21}A_{11}^{-1}(b_1 - A_{12}x_2) + A_{22}x_2 = b_2$$

или после преобразования

$$(A_{22} - A_{21}A_{11}^{-1}A_{12})x_2 = b_2 - A_{21}A_{11}^{-1}b_1. \tag{3.16}$$

В результате система (3.14) преобразовалась к системе

$$\begin{aligned} A_{11}x_1 + A_{12}x_2 &= b_1, \\ (A_{22} - A_{21}A_{11}^{-1}A_{12})x_2 &= b_2 - A_{21}A_{11}^{-1}b_1. \end{aligned} \tag{3.17}$$

(Неизвестные x_1 исключены из второй группы уравнений).

Из (3.17) вроде бы следует, что для реализации блочного исключения нужно вычислять A_{11}^{-1} . На самом деле явно это делать вовсе не обязательно. Принимая во внимание (3.15), введем следующие обозначения:

$$A_{11}^{-1}b_1 = \overset{\circ}{x}_1, \quad A_{11}^{-1}A_{12} = Z_{12}. \tag{3.18}$$

Тогда вторая группа уравнений (3.17) примет вид

$$(A_{22} - A_{21}Z_{12})x_2 = (b_2 - A_{21}\overset{\circ}{x}_1). \tag{3.19}$$

Соотношения (3.18) можно переписать в виде системы уравнений

$$A_{11}\overset{\circ}{x}_1 = b_1, \quad A_{11}Z_{12} = A_{12}, \tag{3.20}$$

а из (3.15) и (3.18) находим, что

$$x_1 = \overset{\circ}{x}_1 - Z_{12}x_2. \tag{3.21}$$

Итак, чтобы найти решение системы (3.14) нужно:

- 1° решить $(m + 1)$ систему (3.20) с матрицей A_{11} для отыскания вектора $\overset{\circ}{x}_1$ и столбцов матрицы Z_{12} ,
- 2° по найденным $\overset{\circ}{x}_1$ и Z_{12} сформировать матрицу и правую часть системы (3.19) и решить полученную систему — найти вектор x_2 ,
- 3° найти вектор x_1 по формулам (3.21).

Замечание 3.3. В трактовке (3.19),(3.20) метода блочного исключения фактически исключенными оказываются не неизвестные x_1 , а неизвестные x_2 . Из системы (3.14) как бы исключается часть неизвестных (именно x_2), затем она решается относительно оставшихся неизвестных (3.20) (но не полностью — нужен еще шаг (3.21)) и лишь потом находится x_2 из (3.19). Отсюда второе название метода — метод частичного исключения неизвестных (исключение x_2).

Пример 3.1. Пусть матрица A имеет портрет, изображенный на рис. 1. Матрица A не является ленточной, хотя ее подматрица A_{11} , расположенная в первых $(n - 1)$ строках и $(n - 1)$ столбцах является ленточной с полушириной $p = 2$. Для решения системы (3.1) с такой матрицей не годится ленточный вариант исключения Гаусса, а применение общего метода требует $O(n^3)$ умножений и делений. Но если можно применить алгоритм блочного исключения,

$$\begin{bmatrix} * & * & * & & & & & & & * \\ * & * & * & * & & & & & & * \\ * & * & * & * & * & & & & & * \\ & & * & * & * & * & * & & & * \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ & & & & & & * & * & * & * \\ * & * & * & * & * & \dots & * & * & * & * \end{bmatrix}$$

Рис. 1

то для решения двух систем (3.20) с пятидиагональными матрицами с использованием ленточного варианта исключения потребуется $O(n)$ действий. Столько же действий потребуется для вычислений по формулам (3.19) и (3.21). В результате система будет решена за $O(n)$ действий.

Пример 3.2. Матрица имеет портрет, изображенный на рис. 2.

$$\begin{bmatrix} * & * & * & * & * & * & \dots & * & * & * \\ * & * & * & & & & & & & * \\ * & * & * & * & & & & & & * \\ * & & * & * & * & & & & & * \\ * & & & * & * & * & & & & * \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ * & & & & & & & * & * & * \\ * & * & * & * & * & * & \dots & * & * & * \end{bmatrix}$$

Рис. 2

Переставляя первую строку на последнее место и то же делая с первым столбцом, получим матрицу с портретом, изображенным на рис. 3.

$$\begin{bmatrix} * & * & & & & & & * & * \\ * & * & * & & & & & * & * \\ & * & * & * & & & & * & * \\ & & * & * & * & & & * & * \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ & & & & & & * & * & * & * \\ * & * & * & * & * & \dots & * & * & * & * \\ * & * & * & * & * & \dots & * & * & * & * \end{bmatrix}$$

Рис. 3

Теперь в качестве A_{11} следует выбрать трехдиагональную матрицу размеров $(n-2) \times (n-2)$, стоящую в левом верхнем углу.

3.4 Обращение матрицы

Сразу же заметим, что обращение матрицы в действительности нужно не так часто, как это может казаться. Предположим, например, что, последующее (за вычислением) использование A^{-1} предусматривает только формирование ее произведений с векторами: $u = A^{-1}r$, $v = A^{-1}s$ и т.д. В таком случае можно было бы вовсе не вычислять A^{-1} в явном виде. Вместо этого достаточно провести для A прямой ход метода Гаусса, а затем находить векторы u , v , ... как решения линейных систем с матрицей A и правыми частями r , s , ... Мы видели в лекции 1, что решение нашей системы требует $q \approx n^2$ операций умножения, но такова же стоимость вычисления произведения A^{-1} на вектор. Предварительная же работа резко различается по объему: $Q \approx n^3/3$ операций умножения для треугольного разложения A и (как мы сейчас покажем) $\approx n^3$ операций для вычисления A^{-1} , т.е. во втором случае втрое больше.

Может случиться, что нужны в явном виде некоторые элементы обратной матрицы, причем они расположены в одном или нескольких (небольшом числе) столбцов A^{-1} . Если номера этих столбцов k , l и т.д., то указанные столбцы опять-таки проще всего вычислить как решения систем с матрицей A , в которых правыми частями служат единичные векторы e_k , e_l и т.д. (столбцы единичной матрицы).

И только если необходимы все или большая часть элементов A^{-1} , применение процедуры численного обращения A оправдана.

Опишем одну из возможных процедур вычисления A^{-1} , основанную на использовании LU -разложения. Пусть $A = LU$. Тогда $A^{-1} = U^{-1}L^{-1}$ или

$$UA^{-1} = L^{-1}.$$

Воспользуемся этим соотношением для вычисления A^{-1} . Допустим сначала, что L^{-1} — известна. Обозначим $A^{-1} = X$, $L^{-1} = Y$. Пусть $x_{.j}$ и $y_{.j}$ — j -е столбцы матриц X и Y , соответственно, т.е. $X = [x_{.1}x_{.2}\dots x_{.n}]$, $x_{.j} = [x_{1j}x_{2j}\dots x_{nj}]^T$. Тогда получим n систем вида

$$Ux_{.j} = y_{.j}, \quad j = 1, \dots, n \quad (3.22)$$

с треугольной матрицей, решения которых могут быть найдены по формулам (1.21) и

$$x_{kj} = \frac{1}{u_{kk}} \left[y_{kj} - \sum_{m=k+1}^n u_{km}x_{mj} \right], \quad k = n, n-1, \dots, 1.$$

Найдем теперь $L^{-1} = Y$. Поскольку $LY = I$, то имеем n систем

$$Ly_{.j} = e_j, \quad j = 1, \dots, n. \quad (3.23)$$

Заметим, что матрица L^{-1} — нижняя треугольная и поэтому у столбца $y_{.j}$ первые $(j-1)$ элемента известны и равны нулю, т.е. $y_{.j}^T = [0 \dots 0 y_{jj} y_{j+1j} \dots y_{nj}] = [0 \bar{y}_{.j}^T]$. Отсюда следует, что для отыскания истинных неизвестных вектора $y_{.j}$ нужно решить систему с треугольной матрицей размеров $(n-j+1) \times (n-j+1)$ относительно $\bar{y}_{.j}$. Для этого можно воспользоваться формулами типа (1.21).

Оценим объем работы по вычислению A^{-1} . В силу (1.22) факторизация $A = LU$ требует $\approx n^3/3$ умножений, решение одной системы (3.22) (см. (1.23)) — $\approx n^2/2$, а всех $\approx n^3/2$, решение всех систем (3.23)

$$\sum_{j=1}^n \frac{(n-j+1)^2}{2} = \frac{1}{2} \sum_{j=1}^n j^2 = \frac{n(n+1)(2n+1)}{12} \approx n^3/6.$$

Складывая, находим, что для вычисления матрицы A^{-1} при помощи описанного алгоритма требуется $\sim n^3$ умножений. Это всего лишь в три раза больше, чем для решения системы (3.1).

§ 4

Устойчивость вычислительных алгоритмов линейной алгебры

4.1 Введение

Исследуем вопрос об устойчивости решения линейной системы по отношению к возмущению правой части. Пусть рассматривается система с квадратной невырожденной матрицей

$$Ax = b \quad (4.1)$$

и система с возмущенной правой частью

$$A\tilde{x} = \tilde{b}. \quad (4.2)$$

Обозначим $\tilde{b} - b = \delta b$, $\tilde{x} - x = \delta x$ и оценим δx через δb . Вычитая (4.1) из (4.2), будем иметь

$$A\delta x = \delta b \quad \Rightarrow \quad \delta x = A^{-1}\delta b. \quad (4.3)$$

Пусть $\|\cdot\|$ — некоторая норма вектора. В линейной алгебре наиболее часто используются следующие нормы

$$\|x\|_\infty = \max_i |x_i|, \quad \|x\|_1 = \sum_{i=1}^n |x_i|, \quad \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}.$$

Как известно, норма матрицы, подчиненная векторной норме $\|\cdot\|$, определяется соотношением

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|. \quad (4.4)$$

Указанным векторным нормам подчинены следующие матричные нормы:

$$\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|, \quad \|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|, \quad \|A\|_2 = \sqrt{\lambda_{\max}(AA^T)}.$$

Очевидно, что для симметричных матриц

$$A = A^T, \quad \|\cdot\|_\infty = \|\cdot\|_1, \quad \text{а } \|A\|_2 = |\lambda_{\max}|,$$

ибо $Ax = \lambda x$, $A^2x = \lambda Ax = \lambda^2x$.

Упражнение 4.1. Доказать, что подчиненные нормы задаются именно этими соотношениями.

Из определения (4.4) матричной нормы, в частности, следует, что

$$\|A\| \geq \frac{\|Ax\|}{\|x\|} \Rightarrow \|Ax\| \leq \|A\| \|x\|. \quad (4.5)$$

Применяя это неравенство ко второму соотношению (4.3), будем иметь

$$\|\delta x\| \leq \|A^{-1}\| \|\delta b\|. \quad (4.6)$$

Соотношение (4.6) дает оценку абсолютной погрешности решения через абсолютную погрешность правой части. При этом множителем (коэффициентом усиления) выступает норма обратной матрицы. Чем больше эта норма, тем на меньшую точность мы можем рассчитывать.

Получим теперь оценку относительной погрешности. Из (4.1) в силу (4.5)

$$\|A\| \|x\| \geq \|b\|. \quad (4.7)$$

Деля (4.6) на (4.7), получим

$$\frac{\|\delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}. \quad (4.8)$$

Это и есть оценка относительной погрешности, которая, вообще говоря, несет больше информации, чем оценка (4.6). Здесь коэффициентом усиления выступает число

$$\|A\| \|A^{-1}\| =: \text{cond } A,$$

называемое числом обусловленности матрицы A .

Если число обусловленности матрицы A большое, то про матрицу A говорят, что она плохо обусловлена. В противном случае говорят о хорошо обусловленной матрице. Поскольку $AA^{-1} = I$, то $\|A\| \|A^{-1}\| \geq 1$, т.е. число обусловленности не может быть меньше единицы. Имея систему с хорошо обусловленной матрицей, мы вправе рассчитывать на то, что при небольших возмущениях правой части возмущение решения не будет слишком велико.

4.2 Примеры систем с плохо обусловленной матрицей

Пример 4.1.

$$\begin{cases} x_1 & = 1, \\ x_1 + 0.01x_2 & = 1. \end{cases} \quad (4.9)$$

Очевидно, что эта система невырождена и ее единственным решением является вектор $[1, 0]^T$. Возмущим правую часть системы (4.9) и найдем решение возмущенной задачи

$$\begin{cases} \tilde{x}_1 & = 1, \\ \tilde{x}_1 + 0.01\tilde{x}_2 & = 1.01, \end{cases} \quad \delta b_2 = 0.01. \quad (4.10)$$

Очевидно, что решением этой системы является вектор $\tilde{x} = [1, 1]^T$, который мало похож на невозмущенный вектор x , ибо

$$\delta x_2 = 1, \quad \delta x_1 = 0, \quad \|\delta x\|_1 = 1.$$

Это значение абсолютной погрешности решения полностью согласуется с оценками (4.6), (4.8), ибо

$$\begin{aligned} A^{-1} &= \begin{bmatrix} 1 & 0 \\ -100 & 100 \end{bmatrix}, \\ \|A\|_1 &= \max_j \sum_{i=1}^2 |a_{ij}| = 2, \\ \|A^{-1}\|_1 &= 101, \quad \|x\|_1 = 1, \quad \|b\|_1 = 2, \quad \|\delta b\|_1 = 0.01 \end{aligned}$$

и, следовательно, в силу (4.6)

$$\|\delta x\|_1 \leq 101 \cdot 0.01 = 1.01,$$

а в силу (4.8)

$$\frac{\|\delta x\|_1}{\|x\|_1} \leq 2 \cdot 101 \cdot \frac{0.01}{2} = 1.01.$$

При заданном (4.10) уровне погрешности правой части обусловленность матрицы этой системы ($\text{cond } A = 202$) следует признать плохой.

Пример 4.2.

$$A = \begin{bmatrix} 1 & -1 & -1 & \dots & -1 & -1 \\ 0 & 1 & -1 & \dots & -1 & -1 \\ 0 & 0 & 1 & \dots & -1 & -1 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 & -1 \\ 0 & 0 & 0 & \dots & 0 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} -1 \\ -1 \\ -1 \\ \dots \\ -1 \\ 1 \end{bmatrix}.$$

В развернутом виде система запишется так

$$\begin{array}{rcccccl}
 x_1 & -x_2 & -x_3 & -\dots & -x_n & = & -1, \\
 & x_2 & -x_3 & -\dots & -x_n & = & -1, \\
 & \dots & \dots & \dots & \dots & & \dots \\
 & & & & x_{n-1} & -x_n & = & -1, \\
 & & & & & x_n & = & 1.
 \end{array} \tag{4.11}$$

Очевидно, что решением системы (4.11) является вектор

$$x = [0, 0, \dots, 0, 1]^T.$$

Легко видеть, что $\det A = 1$.

Возмутим последнюю компоненту вектора b

$$\tilde{b} = [-1, -1, \dots, -1, 1 + \varepsilon]^T$$

и оценим погрешность решения.

Вычитая из возмущенной системы систему (4.11), для погрешности решения получим

$$\begin{array}{rcccccl}
 \delta x_1 & -\delta x_2 & -\dots & -\delta x_n & = & 0, \\
 & \delta x_2 & -\dots & -\delta x_n & = & 0, \\
 & \dots & \dots & \dots & & \dots \\
 & & & \delta x_{n-1} & -\delta x_n & = & 0, \\
 & & & & \delta x_n & = & \varepsilon.
 \end{array}$$

Отсюда находим, что

$$\begin{aligned}
 \delta x_n &= \varepsilon, & \delta x_{n-1} &= \varepsilon, & \delta x_{n-2} &= \delta x_n + \delta x_{n-1} = 2\varepsilon, \\
 \delta x_{n-3} &= \delta x_n + \delta x_{n-1} + \delta x_{n-2} = 4\varepsilon = 2^2\varepsilon.
 \end{aligned}$$

Погрешность в каждой из последующих компонент, начиная с δx_{n-2} , удваивается, так что

$$\delta x_{n-k} = \delta x_n + \delta x_{n-1} + \dots + \delta x_{n-(k-1)} = 2^{k-1}\varepsilon,$$

а

$$\delta x_1 = 2^{n-2}\varepsilon.$$

Таким образом,

$$\begin{aligned}
 \|\delta x\|_\infty &= 2^{n-2}|\varepsilon|, & \|x\|_\infty &= 1, \\
 \|\delta b\|_\infty &= |\varepsilon|, & \|b\|_\infty &= 1, & \|A\|_\infty &= n.
 \end{aligned}$$

Поскольку в силу (4.8)

$$\frac{\|\delta x\|_\infty}{\|x\|_\infty} \leq \text{cond } A \frac{\|\delta b\|_\infty}{\|b\|_\infty},$$

а в рассматриваемом случае

$$\frac{\|\delta x\|_\infty}{\|x\|_\infty} = 2^{n-2}|\varepsilon|,$$

то

$$\text{cond } A = \|A^{-1}\|_{\infty} \|A\|_{\infty} \geq 2^{n-2}$$

и, следовательно,

$$\|A^{-1}\|_{\infty} \geq n^{-1} 2^{n-2}.$$

При $n = 102$, $\|A\|_{\infty} = 102$, $\text{cond } A \geq 2^{100} > 10^{30}$, $\|A^{-1}\| > 10^{27}$. Если $\varepsilon = 10^{-15}$, то $\|\delta x\|_{\infty} > 10^{15}$. Матрица рассматриваемой системы очень плохо обусловлена.

Понятие числа обусловленности введено нами только для невырожденных матриц. Условие $\det A = 0$ означает вырожденность матрицы A , и может сложиться впечатление, что, если $\det A \approx 0$, то матрица плохо обусловлена. Однако, прямой связи между величиной определителя матрицы A и ее обусловленностью нет. Так, определитель матрицы из примера (4.2) равен единице, а

$$\text{cond } A \geq 2^{n-2}.$$

С другой стороны, хорошо обусловленная матрица может иметь очень маленький определитель. Например, у матрицы

$$A = \begin{bmatrix} 10^{-1} & & & \\ & 10^{-1} & & 0 \\ & & \ddots & \\ 0 & & & 10^{-1} \end{bmatrix}$$

$\text{cond } A = 1$, хотя $\det A = 10^{-n}$.

4.3 Пример хорошо обусловленной системы

Рассмотрим следующую систему

$$\begin{aligned} 2x_1 - 9x_2 + 5x_3 &= -4, \\ 1.2x_1 - 5.4x_2 + 6x_3 &= 0.6, \\ x_1 - x_2 - 7.5x_3 &= -8.5. \end{aligned} \tag{4.12}$$

Легко проверить, что решением этой системы является вектор

$$x = [0 \quad 1 \quad 1]^T.$$

Возьмем матрицу и вектор правой части системы (4.12) так, чтобы

$$\delta A = \begin{bmatrix} 0 & 0 & 0 \\ 0 & \varepsilon & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad \delta b = \begin{bmatrix} 0 \\ \varepsilon \\ 0 \end{bmatrix}. \tag{4.13}$$

Очевидно, что решение возмущенной системы совпадает с точным решением, т.е.

$$\tilde{x} = x = [0 \quad 1 \quad 1]^T.$$

Будем теперь решать систему (4.12) с возмущением (4.13) при $\varepsilon = 10^{-4}$ на шестизначном десятичном калькуляторе методом последовательных исключений, используя формулы (1.3) – (1.6).

Прямой ход. 1-й шаг. Вычисляем по формуле (1.5) множители

$$l_{21} = \frac{a_{21}}{a_{11}} = 0.6, \quad l_{31} = \frac{a_{31}}{a_{11}} = 0.5,$$

являющиеся элементами первого столбца левой треугольной матрицы L , и вычитая из второго и третьего уравнений возмущенной системы первое уравнение, умноженное на l_{21} и l_{31} , соответственно, т.е. преобразовывая систему при помощи формул (1.3), (1.4), получим

$$\begin{aligned} 2x_1 - 9x_2 + 5x_3 &= -4, \\ 0.0001x_2 + 3x_3 &= 3.0001, \\ 3.5x_2 - 10x_3 &= -6.5. \end{aligned} \quad (4.14)$$

Все вычисления на этом шаге выполнены точно, без округлений, ибо все числа в промежуточных вычислениях, равно как и окончательные числа, имели мантиссы с меньшим, чем шесть, числом разрядов.

2-й шаг. После вычисления по формуле (1.5) множителя

$$l_{32} = a_{32}^{(1)} / a_{22}^{(1)} = 3.5 / 0.0001 = 35000$$

последнее уравнение системы (4.14) согласно (1.2) должно быть преобразовано к виду

$$a_{33}^{(2)} x_3 = b_3^{(2)}, \quad (4.15)$$

где согласно (1.3)

$$a_{33}^{(2)} = a_{33}^{(1)} - l_{32} a_{23}^{(1)} = -10 - l_{32} \cdot 3 = -10 - 105000 = -105010,$$

а согласно (1.4)

$$\begin{aligned} b_3^{(2)} &= b_3^{(1)} - l_{32} b_2^{(1)} = -6.5 - l_{32} \cdot 3.0001 = \\ &= -6.5 - 105003.5 = -105010, \end{aligned}$$

т.е.

$$-105010 x_3 = -105010.$$

Однако на используемом калькуляторе будет получено уравнение

$$-105010 x_3 = -105011. \quad (4.16)$$

В самом деле, коэффициент $a_{33}^{(2)}$ вычисляется точно, так как при его вычислении не возникает чисел, мантиссы которых имеют больше шести разрядов. В то же время

при вычислении $b_3^{(2)}$ умножение 3.0001 на l_{32} дает семизначное число 105003.5, после округления которого до шести разрядов получим 105004. Вычисление $b_3^{(2)}$ завершается выполнением операции вычитания

$$b_3^{(2)} \approx -6.5 - 105004 = -105010,5 \approx 105011,$$

которая также проводится с округлением, что и приводит к (4.16).

Обратный ход. Из (4.16)

$$\tilde{x}_3 = 1.000009522 \dots \approx 1.00001.$$

Сравнение с истинным значением x_3 показывает хорошую точность. Далее, согласно (1.6), (4.14)

$$x_2 = (3.0001 - 3x_3)/0.0001 = (3.0001 - 3.00003)/0.0001 = 0.7.$$

Здесь все вычисления выполнены точно. Наконец,

$$\begin{aligned} \tilde{x}_1 &= (-4 + 9x_2 - 5x_3)/2 = (-4 + 6.3 - 5.00005)/2 = \\ &= -1.350025 \approx -1.35003. \end{aligned}$$

Итак, приближенное решение $[-1.35003 \quad 0.7 \quad 1.00001]$ мало похоже на точное решение.

В чем причина появления столь значительной погрешности? Говорить о накоплении ошибок округления не приходится, так как всего было выполнено 28 арифметических операций и лишь в четырех случаях потребовалось округление. Предположение о плохой обусловленности системы также не подтверждается, ибо, как показывают вычисления,

$$A = \begin{bmatrix} 2 & -9 & 5 \\ 1.2 & -5.4 & 6 \\ 1 & -1 & -1.5 \end{bmatrix}, \quad \det A = -21,$$

$$A^{-1} = \frac{1}{21} \begin{bmatrix} -46.5 & 72.5 & 27 \\ -15 & 20 & 6 \\ -4.2 & 7 & 0 \end{bmatrix}$$

и, следовательно,

$$\|A\|_1 = 18.5, \quad \|A^{-1}\|_1 = \frac{99.5}{21}, \quad \text{cond } A \leq 10^2.$$

Действительная причина состоит в том, что метод исключения Гаусса в том виде, в каком он был описан, является неустойчивым методом. Чтобы определить, в чем именно его слабость, рассмотрим более внимательно процедуру вычисления $b_3^{(2)}$, где и появились первые округления. При вычислении произведения $l_{32} \cdot b_2^{(1)} = 105003.5$ из-за того, что его мантисса содержит более шести знаков, пришлось прибегнуть к округлению. А так как это произведение к тому же оказалось очень большим, погрешность округления также оказалась больше 0.5.

4.4 Метод Гаусса с выбором ведущего элемента

Из-за отмеченной неустойчивости метод Гаусса в вычислительной практике обычно применяется в сочетании с некоторой схемой выбора ведущего элемента. Например, схема выбора ведущего элемента по столбцу состоит в следующем. Перед началом первого шага среди коэффициентов $a_{11}, a_{21}, \dots, a_{n1}$, образующих первый столбец матрицы A , выбирается коэффициент с наибольшим модулем; пусть это будет $a_{k,1}$. Если $k > 1$, то в системе (4.1) переставляются 1-е и k -е уравнения, при $k = 1$ перестановка не нужна. После этой предварительной работы обычным образом проводится 1-й шаг прямого хода. До начала 2-го шага среди коэффициентов $a_{22}^{(1)}, a_{31}^{(1)}, \dots, a_{n2}^{(1)}$ (т.е. во втором столбце текущей матрицы) выбирается коэффициент $a_{l2}^{(1)}$ с наибольшим модулем. В случае $l > 2$ переставляются 2-е и l -е уравнения, затем выполняется 2-й шаг. И т.д.

Что достигается выбором ведущего элемента по столбцу? Мы можем теперь гарантировать, что множители l_{ij} всех шагов по абсолютной величине ограничены единицей. Формулы (1.3) показывают, что, во-первых, добавки $l_{ik}a_{kj}^{(k-1)}$ к текущим значениям коэффициентов имеют тот же порядок величины, что и сами коэффициенты, во-вторых, за один шаг уровень коэффициентов матрицы может вырасти не более, чем в два раза. Действительно, согласно (1.3)

$$|a_{ij}^{(k)}| = |a_{ij}^{(k-1)} - l_{ik}a_{kj}^{(k-1)}| \leq |a_{ij}^{(k-1)}| + |a_{kj}^{(k-1)}| \leq 2 \max_{ij} |a_{ij}^{(k-1)}|. \quad (4.17)$$

Упражнение 4.2. Решить систему (4.12) методом Гаусса с выбором ведущего элемента по столбцу на 6-разрядном десятичном калькуляторе.

Иногда используется и другая схема выбора ведущего элемента, а именно, схема выбора по строке. Здесь до начала 1-го шага определяется наибольший по модулю среди коэффициентов $a_{11}, a_{12}, \dots, a_{1n}$. Пусть им будет коэффициент a_{1k} . Если $k > 1$, то производится перенумерация неизвестных: 1-е и k -е неизвестные меняются номерами. Это соответствует перестановке столбцов матрицы системы. При $k = 1$ перестановка не нужна. Теперь обычным образом проводится 1-й шаг прямого хода. И т.д.

Переходя к выбору ведущего элемента по столбцу, мы получили для системы (4.12) приближенное решение хорошего качества. Но это не значит, что описанные схемы с выбором ведущего элемента придают методу Гаусса гарантированную устойчивость. Хотя обычно схемы с выбором ведущего элемента по столбцу или по строке действительно обеспечивают устойчивое вычисление.

В каких же случаях утрачивается устойчивость? Чтобы понять это, заметим, что во многих численных методах ошибки промежуточных вычислений в совокупности равносильны тому, как если бы тем же методом точно решали исходную задачу, предварительно изменив ее входные данные. Это относится и к методу Гаусса. Можно показать, что решение линейной системы (4.1), вычисленное методом Гаусса (с той или иной схемой выбора ведущего элемента или вообще без выбора) при наличии ошибок

округления точно удовлетворяет измененному уравнению

$$(A + \delta A)\tilde{x} = b. \quad (4.18)$$

В пояснение сказанного рассмотрим умножение двух треугольных матриц с учетом ошибок округления

$$\begin{aligned} \widetilde{AB} &= \begin{bmatrix} a_{11} & a_{12} \\ 0 & a_{22} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ 0 & b_{22} \end{bmatrix} = \\ &= \begin{bmatrix} a_{11}b_{11}(1 + \varepsilon_1) & (a_{11}b_{12}(1 + \varepsilon_2) + a_{12}b_{22}(1 + \varepsilon_3))(1 + \varepsilon_4) \\ 0 & a_{22}b_{22}(1 + \varepsilon_5) \end{bmatrix} \end{aligned}$$

Если положить

$$\tilde{A} = \begin{bmatrix} a_{11} & a_{12}(1 + \varepsilon_3)(1 + \varepsilon_4) \\ 0 & a_{22}(1 + \varepsilon_5) \end{bmatrix}, \quad \tilde{B} = \begin{bmatrix} b_{11}(1 + \varepsilon_1) & b_{12}(1 + \varepsilon_2)(1 + \varepsilon_4) \\ 0 & b_{22} \end{bmatrix},$$

то легко проверить, что

$$\widetilde{AB} = \tilde{A}\tilde{B}.$$

Для нормы матрицы возмущения в (4.18) справедлива оценка

$$\|\delta A\|_\infty \leq f(n)g(A)\|A\|_\infty p^{-t}. \quad (4.19)$$

Здесь $f(n)$ — некоторая медленно растущая функция от порядка n системы (типа степенной с небольшим показателем); p — основание машинной арифметики; t — число разрядов, отведенных для представления мантиссы. Чтобы определить оставшийся сомножитель $g(A)$, обозначим

$$a = a_0 = \max_{ij} |a_{ij}|, \quad a_k = \max_{ij > k} |a_{ij}^{(k)}|, \quad k = 1, 2, \dots, n-1.$$

Тогда

$$g(A) = \max_{0 \leq k \leq n-1} a_k/a.$$

Таким образом, величина $g(A)$ измеряет насколько могли вырасти элементы промежуточных матриц метода по сравнению с уровнем элементов в исходной матрице A .

Итак, при прочих равных условиях ошибки δA , вносимые методом в исходную информацию, тем больше, чем больший рост элементов матриц допускается в прямом ходе. Если в версии метода Гаусса из первой лекции рост элементов может быть сколь угодно велик, то в схемах с выбором ведущего элемента он ограничен. Действительно, за шаг максимальный модуль элемента матрицы может вырасти самое большее в два раза (см. (4.17)). Так как маловероятно, чтобы возмущение происходило на каждом шаге, то обычно коэффициент роста $g(A)$ невелик. В этом случае благодаря множителю p^{-t} в правой части (4.19) матрицу-возмущение δA можно считать малой по сравнению с A .

На вопрос о том, как сказывается возмущение (4.18) на решении системы (4.1), отвечает

Теорема 4.1. Пусть матрица A системы (4.1) имеет обратную и для ее возмущения δA справедлива оценка

$$\|\delta A\| < \|A^{-1}\|^{-1}. \quad (4.20)$$

Тогда для относительной погрешности решения справедлива оценка

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{cond } A}{1 - \text{cond } A \frac{\|\delta A\|}{\|A\|}} \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right), \quad (4.21)$$

где

$$(A + \delta A)(x + \delta x) = b + \delta b \quad (4.22)$$

— возмущенная система.

Доказательство. Из (4.22)

$$Ax + \delta Ax + A\delta x + \delta A\delta x = b + \delta b.$$

Вычитая отсюда (4.1), получим

$$A\delta x = \delta b - \delta Ax - \delta A\delta x,$$

или

$$\delta x = A^{-1}[\delta b - \delta Ax - \delta A\delta x].$$

Отсюда

$$\|\delta x\| \leq \|A^{-1}\|(\|\delta b\| + \|\delta A\| \|x\| + \|\delta A\| \|\delta x\|).$$

Разрешим это неравенство относительно $\|\delta x\|$

$$(1 - \|A^{-1}\| \|\delta A\|) \|\delta x\| \leq \|A^{-1}\|(\|\delta b\| + \|\delta A\| \|x\|).$$

С учетом (4.20)

$$\|\delta x\| \leq \frac{\|A^{-1}\|(\|\delta b\| + \|\delta A\| \|x\|)}{1 - \|A^{-1}\| \|\delta A\|}. \quad (4.23)$$

Но

$$1 - \|A^{-1}\| \|\delta A\| = 1 - \text{cond } A \frac{\|\delta A\|}{\|A\|}.$$

Учитывая это и деля (4.23) на $\|x\|$, будем иметь

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \|A\| \left(\frac{\|\delta b\|}{\|A\| \|x\|} + \frac{\|\delta A\|}{\|A\|} \right)}{1 - \text{cond } A \frac{\|\delta A\|}{\|A\|}}. \quad (4.24)$$

Поскольку $\|A\| \|x\| \geq \|b\|$, то из (4.24) вытекает (4.21). Теорема доказана.

Приведем примеры матриц, преобразование которых при помощи метода Гаусса с выбором ведущего элемента приводит к максимально возможному увеличению коэффициентов промежуточных матриц прямого хода.

Пример 4.3. Выбор по столбцу

$$A = \begin{bmatrix} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & 1 \end{bmatrix}, \quad A_1 = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & -1 & 1 & 2 \\ 0 & -1 & -1 & 2 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \\ 0 & 0 & -1 & 4 \end{bmatrix},$$

$$A_3 = U = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \\ 0 & 0 & 0 & 8 \end{bmatrix}, \quad L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ -1 & -1 & 1 & 0 \\ -1 & -1 & -1 & 1 \end{bmatrix}.$$

Упражнение 4.3. Показать, что выбор ведущего элемента по строке для матрицы

$$A = \begin{bmatrix} 1 & -1 & -1 & -1 \\ 0 & 1 & -1 & -1 \\ 0 & 0 & 1 & -1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

приводит к максимальному росту элементов промежуточных матриц.

Замечание 4.1. Возможна еще одна схема выбора ведущего элемента: выбор по всей матрице. В этом случае гарантируется полная устойчивость метода Гаусса, однако сама процедура выбора такого ведущего элемента очень трудоемка — для ее реализации требуется $O(n^3)$ действий, что сравнимо с трудоемкостью самого метода и, следовательно, существенно удорожает решение. Отметим, что при выборе ведущего элемента по столбцу или по строке требуется лишь $O(n^2)$ дополнительных операций.

§ 5

Методы вращений и отражений

5.1 Метод вращений

Вновь обратимся к решению системы

$$Ax = b \quad (5.1)$$

с невырожденной матрицей. Рассмотрим метод вращений решения системы (5.1), который позволяет получить представление матрицы A в виде произведения ортогональной матрицы Q и верхней треугольной матрицы R .

Как и в методе Гаусса, в методе вращений на первом шаге неизвестное x_1 исключается из всех уравнений кроме первого. Для того, чтобы исключить x_1 из второго уравнения, умножим первое уравнение на некоторое число c_{12} , а второе — на s_{12} и заменим первое уравнение суммой вновь полученных уравнений. Затем умножим первое уравнение на s_{12} , второе — на c_{12} и вычтем первое из второго:

$$\begin{aligned} (c_{12}a_{11} + s_{12}a_{21})x_1 &+ (c_{12}a_{12} + s_{12}a_{22})x_2 &+ \dots &= c_{12}b_1 + s_{12}b_2, \\ (-s_{12}a_{11} + c_{12}a_{21})x_1 &+ (-s_{12}a_{12} + c_{12}a_{22})x_2 &+ \dots &= -s_{12}b_1 + c_{12}b_2. \end{aligned}$$

Числа c_{12} и s_{12} выберем из условий

$$c_{12}^2 + s_{12}^2 = 1, \quad -s_{12}a_{11} + c_{12}a_{21} = 0. \quad (5.2)$$

Решение нелинейной системы (5.2) при $a_{11}^2 + a_{21}^2 \neq 0$ дают, например, формулы

$$c_{12} = \frac{a_{11}}{\sqrt{a_{11}^2 + a_{21}^2}}, \quad s_{12} = \frac{a_{21}}{\sqrt{a_{11}^2 + a_{21}^2}}. \quad (5.3)$$

Если же $a_{11}^2 + a_{21}^2 = 0$, то решение (5.2) дают числа $c_{12} = 1$, $s_{12} = 0$. Второе из уравнений (5.2) как раз и означает, что в новом втором уравнении неизвестного x_1 не будет.

Это преобразование эквивалентно умножению слева матрицы системы (5.4) и ее вектора правой части на матрицу

$$T_{13} = \begin{bmatrix} c_{13} & 0 & s_{13} & & & \\ 0 & 1 & 0 & & & 0 \\ -s_{13} & 0 & c_{13} & & & \\ & & & 1 & & \\ & 0 & & & \ddots & \\ & & & & & 1 \end{bmatrix}$$

и приводит к тому, что коэффициент при x_1 в преобразованном третьем уравнении обращается в нуль.

Продолжая подобным образом, мы исключим x_1 из всех остальных уравнений. Первое уравнение изменяется на каждом таком "малом" шаге, которых будет $n - 1$. Поэтому, по завершении первого шага система (5.1) примет вид

$$\begin{aligned} a_{11}^{(n-1)} x_1 + a_{12}^{(n-1)} x_2 + a_{13}^{(n-1)} x_3 + \dots + a_{1n}^{(n-1)} x_n &= b_1^{(n-1)}, \\ a_{22}^{(1)} x_2 + a_{23}^{(1)} x_3 + \dots + a_{2n}^{(1)} x_n &= b_2^{(1)}, \\ \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \\ a_{n2}^{(1)} x_2 + a_{n3}^{(1)} x_3 + \dots + a_{nn}^{(1)} x_n &= b_n^{(1)}. \end{aligned}$$

Упражнение 5.1. Вывести соотношения для $a_{1j}^{(l-1)}$, $a_{lj}^{(1)}$, $b_1^{(l-1)}$, $b_l^{(1)}$ и преобразующих коэффициентов c_{1l} и s_{1l} , $l = 2, \dots, n$; $j = 1, \dots, n$.

Ответ.

$$\begin{aligned} a_{1j}^{(l-1)} &= c_{1l} a_{1j}^{(l-2)} + s_{1l} a_{lj}, & a_{lj}^{(1)} &= -s_{1l} a_{1j}^{(l-2)} + c_{1l} a_{lj}, \\ & & j &= 1, 2, \dots, n, \\ b_1^{(l-1)} &= c_{1l} b_1^{(l-2)} + s_{1l} b_l, & b_l^{(1)} &= -s_{1l} b_1^{(l-2)} + c_{1l} b_l, \\ & & l &= 2, 3, \dots, n, \\ c_{1l} &= \frac{a_{11}^{(l-2)}}{\sqrt{(a_{11}^{(l-2)})^2 + a_{l1}^2}}, & s_{1l} &= \frac{a_{l1}}{\sqrt{(a_{11}^{(l-2)})^2 + a_{l1}^2}}, \\ & & l &= 2, 3, \dots, n, \quad a_{11}^{(0)} = a_{11}. \end{aligned}$$

Замечание 5.1. Поскольку в силу невырожденности матрицы A по крайней мере один из коэффициентов $a_{i1} \neq 0$, то $a_{11}^{(n-1)} \neq 0$.

В матричной записи эта система имеет вид

$$A^{(1)} x = b^{(1)},$$

где

$$A^{(1)} = T_1 A, \quad b^{(1)} = T_1 b, \quad T_1 = T_{1n} T_{1n-1} \dots T_{13} T_{12}.$$

На втором шаге метода вращений из третьего, четвертого и т.д. уравнений полученной системы исключается неизвестное x_2 . Шаг состоит из $(n - 2)$ "малых" шагов, и в каждом из них второе уравнение комбинируется с одним из нижележащих. После выполнения второго шага система преобразуется к виду

$$\begin{aligned} a_{11}^{(n-1)} x_1 + a_{12}^{(n-1)} x_2 + a_{13}^{(n-1)} x_3 + \dots + a_{1n}^{(n-1)} x_n &= b_1^{(n-1)}, \\ a_{22}^{(n-1)} x_2 + a_{23}^{(n-2)} x_3 + \dots + a_{2n}^{(n-2)} x_n &= b_2^{(n-1)}, \\ a_{33}^{(2)} x_3 + \dots + a_{3n}^{(2)} x_n &= b_3^{(2)}, \\ &\dots\dots\dots \\ a_{n3}^{(2)} x_3 + \dots + a_{nn}^{(2)} x_n &= b_n^{(2)}. \end{aligned}$$

Упражнение 5.2. Вывести соотношения для $a_{2j}^{(l-1)}$, $a_{lj}^{(2)}$, $b_2^{(l-1)}$, $b_l^{(2)}$, c_{2l} и s_{2l} .

Ответ.

$$\begin{aligned} a_{2j}^{(l-1)} &= c_{2l} a_{2j}^{(l-2)} + s_{2l} a_{lj}^{(1)}, & a_{lj}^{(2)} &= -s_{2l} a_{2j}^{(l-2)} + c_{2l} a_{2j}^{(1)}, \\ & & & j = 2, 3, \dots, n, \\ b_2^{(l-1)} &= c_{2l} b_2^{(l-2)} + s_{2l} b_l^{(1)}, & b_l^{(2)} &= -s_{2l} b_2^{(l-2)} + c_{2l} b_l^{(1)}, \\ & & & l = 3, 4, \dots, n, \\ c_{2l} &= \frac{a_{22}^{(l-2)}}{\sqrt{\left(a_{22}^{(l-2)}\right)^2 + \left(a_{l2}^{(1)}\right)^2}}, & s_{2l} &= \frac{a_{l2}^{(1)}}{\sqrt{\left(a_{22}^{(l-2)}\right)^2 + \left(a_{l2}^{(1)}\right)^2}}. \end{aligned}$$

В матричной форме эта система имеет вид

$$A^{(2)} x = b^{(2)},$$

где

$$A^{(2)} = T_2 A^{(1)}, \quad b^{(2)} = T_2 b^{(1)}, \quad T_2 = T_{2n} T_{2(n-1)} \dots T_{24} T_{23}.$$

После $(n - 1)$ шагов получим систему

$$\begin{aligned} a_{11}^{(n-1)} x_1 + a_{12}^{(n-1)} x_2 + a_{13}^{(n-1)} x_3 + \dots + a_{1n}^{(n-1)} x_n &= b_1^{(n-1)}, \\ a_{22}^{(n-1)} x_2 + a_{23}^{(n-2)} x_3 + \dots + a_{2n}^{(n-2)} x_n &= b_2^{(n-1)}, \\ &\dots\dots\dots \\ a_{nn}^{(n-1)} x_n &= b_n^{(n-1)}, \end{aligned} \tag{5.7}$$

где

$$\begin{aligned} a_{kj}^{(l-1)} &= c_{kl}a_{kj}^{(l-2)} + s_{kl}a_{lj}^{(k-1)}, & a_{lj}^{(k)} &= -s_{kl}a_{kj}^{(l-2)} + c_{kl}a_{lj}^{(k-1)}, \\ & & j &= k, k+1, \dots, n, \\ b_k^{(l-1)} &= c_{kl}b_k^{(l-2)} + s_{kl}b_l^{(k-1)}, & b_l^{(k)} &= -s_{kl}b_k^{(l-2)} + c_{kl}b_l^{(k-1)}, \\ & & k &= 1, \dots, n, \\ & & l &= k+1, \dots, n, \end{aligned} \quad (5.8)$$

а

$$c_{kl} = \frac{a_{kk}^{(l-2)}}{\sqrt{\left(a_{kk}^{(l-2)}\right)^2 + \left(a_{lk}^{(k-1)}\right)^2}}, \quad s_{kl} = \frac{a_{lk}^{(k-1)}}{\sqrt{\left(a_{kk}^{(l-2)}\right)^2 + \left(a_{lk}^{(k-1)}\right)^2}}. \quad (5.9)$$

В матричной записи полученная система имеет вид

$$A^{(n-1)}x = b^{(n-1)},$$

где

$$A^{(n-1)} = T_{n-1}A^{(n-2)}, \quad b^{(n-1)} = T_{n-1}b^{(n-2)}, \quad T_{n-1} = T_{n-1n}.$$

Обозначим через R полученную верхнюю треугольную матрицу $A^{(n-1)}$. Она связана с исходной матрицей A равенством

$$R = TA, \quad (5.10)$$

где $T = T_{n-1}T_{n-2} \dots T_1$.

Замечание 5.2. Действие матрицы T_{kl} на вектор x эквивалентно его повороту против хода часовой стрелки в координатной плоскости Ox_kx_l на угол φ_{kl} такой, что

$$\cos \varphi_{kl} = c_{kl}, \quad \sin \varphi_{kl} = s_{kl}.$$

Существование такого угла гарантируется соотношениями (5.9). Эта геометрическая интерпретация и дала название методу.

С учетом (5.2), (5.6) и т.д. легко видеть, что

$$T_{kl}T_{kl}^T = I,$$

т.е. матрицы T_{kl} ортогональные. Произведение ортогональных матриц есть матрица ортогональная, и поэтому T есть ортогональная матрица, равно как и $T^T = T^{-1} = Q$. Отсюда и из (5.10) находим, что

$$A = QR. \quad (5.11)$$

Упражнение 5.3. Показать, что для построения разложения (5.11) с использованием формул (5.8), (5.9), требуется $\approx 4n^3/3$ действий умножения.

Замечание 5.3. Принимая во внимание результаты вычислений из упражнения 5.3, заключаем, что метод вращений примерно в четыре раза более трудоемок, чем метод Гаусса.

Выясним, какими же достоинствами обладает этот метод по сравнению с методом Гаусса, благодаря которым он заслуживает право на существование несмотря на большую трудоемкость.

С учетом (5.2) из (5.5) находим, что

$$\begin{aligned} [a_{1j}^{(1)}]^2 + [a_{2j}^{(1)}]^2 &= c_{12}^2 a_{1j}^2 + s_{12}^2 a_{2j}^2 + 2c_{12}s_{12}a_{1j}a_{2j} + \\ &+ s_{12}^2 a_{1j}^2 + c_{12}^2 a_{2j}^2 - 2s_{12}c_{12}a_{1j}a_{2j} = a_{1j}^2 + a_{2j}^2, \\ & j = 1, 2, \dots, n, \end{aligned}$$

т.е. сумма квадратов первых элементов j -го столбца не изменилась. Если учесть, что на первом малом шаге коэффициенты остальных уравнений не изменились, то полученное равенство означает: длина любого столбца матрицы не изменилась. Точно так же из формул (5.6), (5.7) выводим, что

$$[a_{1j}^{(2)}]^2 + [a_{3j}^{(1)}]^2 = [a_{1j}^{(1)}]^2 + a_{3j}^2, \quad j = 1, \dots, n,$$

т.е. длины столбцов неизменны и на 2-ом малом шаге. Это верно и по отношению к каждому из последующих шагов, и по отношению к прямому ходу в целом. Таким образом, в отличие от метода Гаусса метод вращений застрахован от роста элементов промежуточных матриц: на протяжении всего процесса исключения абсолютная величина коэффициента $a_{ij}^{(k)}$ не может превосходить длину j -го столбца исходной матрицы A . Как следствие, матрица δA — матрица эквивалентных возмущений, учитывающая ошибки промежуточных возмущений, для метода вращений будет мала, т.е. метод всегда устойчив.

5.2 Метод отражений

Рассмотрим еще один метод, который дает разложение матрицы A в виде произведения ортогональной и верхней треугольной матриц. Это будет метод отражений.

Пусть w — некоторый вектор (столбец) единичной длины

$$\|w\|_2^2 = (w, w) = w^T w = 1. \quad (5.12)$$

Введем в рассмотрение матрицу

$$U = I - 2ww^T, \quad (5.13)$$

которую назовем матрицей отражения и изучим ее свойства.

1°. Матрица U симметрична, т.е.

$$U = U^T. \quad (5.14)$$

В самом деле, так как

$$(ww^T)^T = (w^T)^T w^T = ww^T,$$

то ww^T есть симметричная матрица, а в силу (5.13) вместе с ней симметричной является и матрица U .

2°. Матрица U есть ортогональная матрица, т.е.

$$U^{-1} = U^T. \quad (5.15)$$

Принимая во внимание (5.14), (5.13) и (5.12), имеем

$$\begin{aligned} UU^T &= UU = (I - 2ww^T)(I - 2ww^T) = \\ &= I - 4ww^T + 4w \underbrace{w^T w}_{\parallel 1} w^T = I, \end{aligned}$$

что и означает справедливость (5.15).

3°. Число $\lambda = -1$ является однократным собственным значением матрицы U , которому отвечает собственный вектор w из (5.13). Число $\lambda = 1$ является $(n-1)$ -кратным собственным значением матрицы U , которому отвечает $(n-1)$ -мерное собственное подпространство, состоящее из всех векторов v , ортогональных w .

В силу (5.14), (5.15) имеем $U^2 = UU^T = I$. Так как все собственные значения матрицы I равны 1, то собственные значения матрицы U удовлетворяют соотношению $\lambda_U^2 = 1$, и, следовательно, равны либо +1, либо -1.

Далее, принимая во внимание (5.13), (5.12), находим, что

$$Uw = (I - 2ww^T)w = w - 2w \underbrace{w^T w}_{\parallel 1} = -w. \quad (5.16)$$

Наконец, пусть

$$(v, w) = w^T v = 0.$$

Тогда

$$Uv = (I - 2ww^T)v = v - 2ww^T v = v, \quad (5.17)$$

что и требовалось доказать.

4°. Вектор Uy есть зеркальное отражение вектора y относительно плоскости, ортогональной вектору w .

Пусть y — произвольный вектор. Представим его в виде

$$y = z + v, \quad (5.18)$$

где z — проекция y на w , т.е. $z = (w, y)w$, а вектор v ортогонален w : $(w, v) = 0$. Принимая во внимание (5.16), (5.17), находим, что

$$Uy = U(z + v) = -z + v \quad (5.19)$$

(см. рис.1)

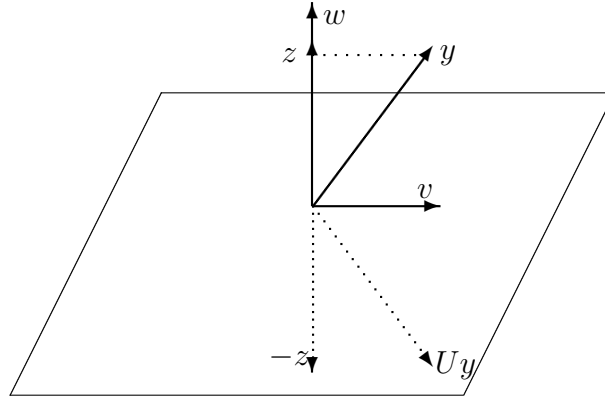


Рис. 1

5°. Векторы $y - Uy$ и $y + Uy$ ортогональны друг другу. При этом, если $y = z + v$, как в (5.18), то

$$y - Uy = 2z, \quad z \parallel w, \quad (5.20)$$

$$y + Uy = 2v, \quad v \perp w. \quad (5.21)$$

Утверждения следуют из (5.18), (5.19).

6°. Пусть x и y — произвольные векторы единичной длины. Тогда, если

$$w = \frac{y + \operatorname{sign}(x, y)x}{\|y + \operatorname{sign}(x, y)x\|_2}, \quad (5.22)$$

то матрица отражения U , построенная по вектору w , переводит вектор y в вектор, коллинеарный вектору x .

В самом деле, пусть матрица отражения U обладает искомым свойством, т.е. $Uy = -\sigma x$. Найдем вектор w , образующий эту матрицу. Согласно свойству 5° (см. (5.20)) искомым вектор w коллинеарен вектору $y - Uy$, а поскольку $\|w\|_2 = 1$, то

$$w = \frac{y + \sigma x}{\|y + \sigma x\|_2}. \quad (5.23)$$

В силу 2° матрица U ортогональна и поэтому

$$\|Uy\|_2 = \|y\|_2 = 1 = |\sigma| \|x\|_2 = |\sigma|,$$

т.е. $\sigma = \pm 1$. Если $y = \pm x$, то знаменатель (5.23) будет равен нулю при $\sigma = -\text{sign}(x, y) = \mp 1$. Этого не произойдет, если $\sigma = \text{sign}(x, y)$. Если вектор y близок к x , то мы будем избавлены от неприятностей, связанных с делением на малое число, выбрав и здесь $\sigma = \text{sign}(x, y)$. Из сказанного следует, что в общем случае целесообразно выбирать σ согласно (5.22), а если $(y, x) = 0$, то, например, $\sigma = 1$.

Воспользуемся матрицей отражения для приведения квадратной матрицы к треугольному виду. На первом шаге приведения рассмотрим в качестве вектора y из свойства 6° нормированный первый столбец матрицы A

$$y_1 = [a_{11} \ a_{21} \ \dots \ a_{n1}]^T / \sqrt{\sum_{i=1}^n a_{i1}^2}, \quad (5.24)$$

а в качестве x — вектор $e_1 = [1 \ 0 \ \dots \ 0]^T$. Если $a_{21} = a_{31} = \dots = a_{n1} = 0$, то переходим к следующему шагу, положив $A^{(1)} = A$, $U_1 = I$ и введя обозначения $a_{ij}^{(1)} = a_{ij}$. В противном случае умножим матрицу A слева на матрицу отражения

$$U_1 = I - 2w_1 w_1^T = I_n - 2w_1 w_1^T, \quad (5.25)$$

где вектор w_1 вычисляется согласно формуле (5.22)

$$w_1 = \frac{y_1 + \text{sign}(e_1, y_1)e_1}{\|y_1 + \text{sign}(e_1, y_1)e_1\|_2}. \quad (5.26)$$

В результате получим матрицу

$$A^{(1)} = U_1 A,$$

в первом столбце которой стоят нули во всех позициях, кроме первой. Этим заканчивается первый этап.

Пусть мы уже осуществили $l-1 > 0$ шагов и пришли к матрице $A^{(l-1)}$ с элементами $a_{ij}^{(l-1)}$ такими, что $a_{ij}^{(l-1)} = 0$ при $i > j$, $j = 1, \dots, l-1$. В пространстве \mathbb{R}_{n-l+1} векторов размерности $n-l+1$ рассмотрим вектор

$$y_l = [a_{ll}^{(l-1)} \ a_{l+1,l}^{(l-1)} \ \dots \ a_{nl}^{(l-1)}]^T / \sqrt{(a_{ll}^{(l-1)})^2 + \dots + (a_{nl}^{(l-1)})^2}.$$

Если $a_{l+1,l}^{(l-1)} = a_{l+2,l}^{(l-1)} = \dots = a_{nl}^{(l-1)} = 0$, то переходим к следующему шагу, положив

$$A^{(l)} = A^{(l-1)}, \quad U_l = I.$$

В противном случае строим матрицу отражения

$$V_l = I_{n-l+1} - 2w_l w_l^T$$

(размеры матрицы V_l и вектора w_l равны $(n-l+1)$), переводящую вектор y_l в вектор, коллинеарный $e_l = [1 \ 0 \ \dots \ 0]^T \in \mathbb{R}_{n-l+1}$, и переходим к матрице

$$A^{(l)} = U_l A^{(l-1)}, \quad (5.27)$$

где

$$U_l = \begin{bmatrix} I_{l-1} & 0 \\ 0 & V_l \end{bmatrix}$$

После $(n - 1)$ шагов мы приходим к матрице

$$A^{(n-1)} = U_{n-1}U_{n-2}\dots U_1A,$$

имеющей правую треугольную форму. Обозначим

$$U_{n-1}\dots U_1 = U.$$

Тогда

$$A^{(n-1)} = UA, \quad A = U^T A^{(n-1)}.$$

Если нужно решить систему (5.1), то после описанных преобразований приходим к эквивалентной системе

$$A^{(n-1)}x = Ub \tag{5.28}$$

с треугольной матрицей.

Упражнение 5.4. Построить алгоритм метода отражений, при котором для разложение матрицы в произведение ортогональной и верхней треугольной матриц достаточно $\approx \frac{2}{3}n^3$ действий умножения.

§ 6

Разностные уравнения

6.1 Разностные уравнения

Пусть $y(n) = y_n$ — функция целочисленного аргумента $n \in \mathbb{Z}$. Будем ее называть сеточной функцией. Обозначим через ∇ (набла) оператор левой конечной разности (разности назад), т.е.

$$\nabla y_n = y_n - y_{n-1}. \quad (6.1)$$

Степень оператора ∇ определим рекуррентным образом

$$\nabla^k = \nabla(\nabla^{k-1}). \quad (6.2)$$

Пусть $a(n)$, $b(n)$, $c(n)$, $d(n)$ и $f(n)$ — заданные сеточные функции. Рассмотрим уравнение

$$a(n)\nabla^3 y_n + b(n)\nabla^2 y_n + c(n)\nabla y_n + d(n)y_n = f(n) \quad (6.3)$$

относительно сеточной функции y_n . Уравнение (6.3) называется разностным уравнением. Разностные уравнения являются аналогами дифференциальных уравнений и в значительной степени повторяют свойства последних. Как и в случае дифференциальных уравнений, важным является понятие порядка разностного уравнения. Если $a(n) \neq 0$, то казалось бы естественным объявить порядком уравнения (6.3) число три. Однако при таком определении порядка разностного уравнения нас ждут неприятности. Чтобы убедиться в этом, положим в (6.3) $a(n) = 1$, $b(n) = 0$, $c(n) = -3$, $d(n) = 2$. В результате получим уравнение

$$\nabla^3 y_n - 3\nabla y_n + 2y_n = f(n). \quad (6.4)$$

Принимая во внимание (6.1) и (6.2), находим, что

$$\nabla y_n = y_n - y_{n-1}, \quad \nabla^3 y_n = y_n - 3y_{n-1} + 3y_{n-2} - y_{n-3},$$

а, подставляя эти выражения в (6.4), будем иметь

$$y_n - 3y_{n-1} + 3y_{n-2} - y_{n-3} - 3y_n + 3y_{n-1} + 2y_n = f(n)$$

или

$$3y_{n-2} - y_{n-3} = f(n). \quad (6.5)$$

Вводя новый индекс $m = n - 2$, уравнение (6.5) преобразуем к виду

$$3y_m - y_{m-1} = f(m + 2). \quad (6.6)$$

Это уравнение эквивалентно уравнению (6.4), и назвать его разностным уравнением третьего порядка просто не поворачивается язык. И дело, конечно, не просто в названии. От удачно введенного определения зависит простота последующих утверждений, использующих это определение. Поскольку запись (6.3) не содержит явным образом информации о числе, которым следовало бы определить порядок разностного уравнения, то будем разностное уравнение записывать в виде

$$\Phi(n, y_n, y_{n-1}, \dots, y_{n-k}) = 0. \quad (6.7)$$

Определение 6.1. Уравнение (6.7) называется разностным уравнением.

Определение 6.2. Разностное уравнение (6.7), если оно явно зависит от y_n и от y_{n-k} , называется уравнением k -го порядка.

Определение 6.3. Разностное уравнение k -го порядка называется линейным, если оно линейно зависит от $y_n, y_{n-1}, \dots, y_{n-k}$.

Мы будем изучать только линейные разностные уравнения, которые будем записывать в виде

$$\sum_{j=0}^k \alpha_j(n) y_{n-j} = f(n), \quad n \in \mathbb{Z}. \quad (6.8)$$

Пока мы предполагаем, что уравнение (6.8) задано при всех $n \in \mathbb{Z}$. Уравнение (6.8) будет уравнением k -го порядка, если коэффициенты $\alpha_0(n)$ и $\alpha_k(n)$ не обращаются в нуль ни при одном $n \in \mathbb{Z}$.

Определение 6.4. Сеточная функция $y_n, n \in \mathbb{Z}$ называется решением уравнения (6.8), если при подстановке ее в (6.8) последнее превращается в тождество.

Определение 6.5. Сеточная функция $y_n, n \in \mathbb{Z}$ называется общим решением разностного уравнения (6.8), если в ней содержится любое решение (6.8).

Для того, чтобы определить какое-либо решение уравнения (6.8) (частное решение) достаточно указать его значения в любых k последовательных точках, например, $n_0, n_0 + 1, \dots, n_0 + k - 1$.

6.2 Линейные разностные уравнения первого порядка

Эти уравнения имеют вид

$$\alpha_0(n)y_n + \alpha_1(n)y_{n-1} = f(n). \quad (6.9)$$

Поскольку $\alpha_0(n) \neq 0$, то на этот коэффициент уравнение можно поделить. Пусть $\alpha_1(n)/\alpha_0(n) = -q_n$, а $f(n)/\alpha_0(n)$ снова обозначим через $f(n) = f_n$. Тогда разностное уравнение первого порядка (6.9) можно переписать так

$$y_n = q_n y_{n-1} + f_n. \quad (6.10)$$

Разрешить разностное уравнение — значит выразить y_n через известные величины. Чтобы можно было решить (6.10), нужно задать начальное условие

$$y_0 = a. \quad (6.11)$$

Используя теперь рекуррентные соотношения (6.10), можно последовательно определить y_n при всех последующих значениях n :

$$\begin{aligned} y_1 &= q_1 y_0 + f_1 = q_1 a + f_1, \\ y_2 &= q_2 y_1 + f_2 = q_2 (q_1 a + f_1) + f_2 = q_1 q_2 a + q_2 f_1 + f_2 \end{aligned}$$

и т.д.

Часто бывает полезно иметь не рекуррентное соотношение для последовательного вычисления решения, а некоторую формулу, представляющую решение. Найдем представление решения уравнения (6.10). Для этого рассмотрим сначала отвечающее ему однородное уравнение

$$y_n = q_n y_{n-1} \quad (6.12)$$

и найдем его решение. Имеем

$$\begin{aligned} y_1 &= q_1 y_0, \\ y_2 &= q_2 y_1, \\ &\dots\dots\dots \\ y_n &= q_n y_{n-1}. \end{aligned}$$

Перемножая последовательно полученные равенства и сокращая левую и правую части на $y_1 y_2 \dots y_{n-1}$, получим

$$y_n = q_1 \dots q_n y_0 = y_0 \prod_{j=1}^n q_j. \quad (6.13)$$

Величина y_0 есть начальное значение y_n и является произвольной постоянной. Решение однородного уравнения (6.12) найдено.

Замечание 6.1. Напомним, что если линейное однородное дифференциальное уравнение первого порядка записать в виде $y' = P(x)y$, то его общее решение примет вид

$$y(x) = c \exp \left\{ \int_0^x P(\xi) d\xi \right\}.$$

Обратимся теперь к неоднородному уравнению (6.10). Его решение будем искать, используя решение однородного уравнения (6.12), методом вариации постоянной. Пусть

$$\overset{\circ}{y}_n = \prod_{j=1}^n q_j. \quad (6.14)$$

Это — решение уравнения (6.12), а $c \overset{\circ}{y}_n$ — его общее решение. Заставим коэффициент c зависеть от n и в таком виде будем искать решение уравнения (6.10)

$$y_n = c_n \overset{\circ}{y}_n. \quad (6.15)$$

Подставляя (6.15) в (6.10), получим

$$c_n \overset{\circ}{y}_n = q_n c_{n-1} \overset{\circ}{y}_{n-1} + f_n.$$

Из (6.12) $\overset{\circ}{y}_n = q_n \overset{\circ}{y}_{n-1}$ и поэтому

$$c_n \overset{\circ}{y}_n = c_{n-1} \overset{\circ}{y}_n + f_n,$$

т.е.

$$c_n = c_{n-1} + f_n / \overset{\circ}{y}_n.$$

Отсюда

$$\begin{aligned} c_1 &= c_0 + f_1 / \overset{\circ}{y}_1, \\ c_2 &= c_1 + f_2 / \overset{\circ}{y}_2, \\ &\dots\dots\dots \\ c_n &= c_{n-1} + f_n / \overset{\circ}{y}_n, \end{aligned}$$

Складывая эти соотношения, находим, что

$$c_n = \sum_{k=1}^n \frac{f_k}{\overset{\circ}{y}_k} + c_0,$$

а принимая во внимание (6.14), будем иметь

$$c_n = \sum_{k=1}^n f_k \prod_{j=1}^k q_j^{-1} + c_0.$$

6.3. УРАВНЕНИЯ K -ГО ПОРЯДКА С ПОСТОЯННЫМИ КОЭФФИЦИЕНТАМИ 61

Подставляя это выражение в (6.15), получим общее решение неоднородного уравнения (6.10)

$$y_n = \prod_{j=1}^n q_j \left(c + \sum_{k=1}^n f_k \prod_{j=1}^k q_j^{-1} \right). \quad (6.16)$$

Замечание 6.2. Напомним, что если линейное неоднородное дифференциальное уравнение первого порядка записать в виде $y' = P(x)y + f(x)$, то его общее решение примет вид

$$y(x) = \exp \left\{ \int_0^x P(\xi) d\xi \right\} \left(c + \int_0^x f(\eta) \exp \left\{ - \int_0^\eta P(\xi) d\xi \right\} d\eta \right).$$

Если коэффициент $q_n = \text{const} = q$, то из (6.16) находим, что

$$y_n = q^n \left(c + \sum_{k=1}^n f_k q^{-k} \right), \quad (6.17)$$

а если и $f_n = \text{const} = f$, то при $q \neq 1$

$$y_n = q^n \left(c + f \sum_{k=1}^n q^{-k} \right) = q^n \left(c + f \frac{q^{-1} - q^{-n+1}}{1 - q^{-1}} \right) = cq^n + f \frac{1 - q^n}{1 - q}. \quad (6.18)$$

6.3 Разностные уравнения k -го порядка с постоянными коэффициентами

Если коэффициенты $\alpha_j(n)$ из (6.8) не зависят от n , то мы имеем разностное уравнение с постоянными коэффициентами

$$\sum_{j=0}^k \alpha_j y_{n-j} = f(n), \quad n \in \mathbb{Z}. \quad (6.19)$$

Решение отвечающего (6.19) однородного уравнения

$$\sum_{j=0}^k \alpha_j y_{n-j} = 0 \quad (6.20)$$

можно искать в виде

$$y_n = q^n, \quad q \neq 0 \quad (\text{ср. с } y(x) = e^{\lambda x}), \quad (6.21)$$

где $q = \text{const} \neq 0$. Подставляя (6.21) в (6.20), получим

$$q^{n-k} \sum_{j=0}^k \alpha_j q^{k-j} = 0.$$

На q^{n-k} можно сократить, в результате чего для отыскания q получим алгебраическое уравнение степени k

$$\alpha_0 q^k + \alpha_1 q^{k-1} + \dots + \alpha_{k-1} q + \alpha_k = 0, \quad (6.22)$$

называемое характеристическим уравнением, отвечающим разностному уравнению (6.20).

Характеристическое уравнение (6.22) имеет ровно k корней, включая кратные и комплексные. Обозначим их через

$$q_1, q_2, \dots, q_k. \quad (6.23)$$

Очевидно, что сеточные функции

$$q_l^n, \quad l = 1, \dots, k \quad (6.24)$$

являются решениями разностного уравнения (6.20).

Имеет место

Теорема 6.1. *Если корни (6.23) характеристического уравнения (6.22) простые, то решения (6.24) разностного уравнения (6.20) линейно независимы, а общее решение этого уравнения имеет вид*

$$y_n = \sum_{l=1}^k c_l q_l^n.$$

Доказательство. Проведем доказательство линейной независимости (6.20) при $k = 2$. Допустим противное, т.е. пусть

$$c_1 q_1^n + c_2 q_2^n \equiv 0, \quad |c_1| + |c_2| \neq 0.$$

Но тогда и

$$c_1 q_1^{n-1} + c_2 q_2^{n-1} = 0.$$

Рассмотрим эти два тождества как систему уравнений для определения c_1 и c_2 . Находим, что определитель этой системы

$$\Delta = \begin{vmatrix} q_1^n & q_2^n \\ q_1^{n-1} & q_2^{n-1} \end{vmatrix} = (q_1 q_2)^{n-1} (q_1 - q_2) \neq 0$$

и, следовательно, система имеет лишь тривиальное решение $c_1 = c_2 = 0$. Это противоречит предположению, что и доказывает теорему.

Замечание 6.3. Если комплексное число $q = |q|e^{i\varphi}$, $\varphi \neq m\pi$, $m \in \mathbb{Z}$ является корнем характеристического уравнения (6.22), коэффициенты которого действительны, то

6.3. УРАВНЕНИЯ К-ГО ПОРЯДКА С ПОСТОЯННЫМИ КОЭФФИЦИЕНТАМИ 63

число $\bar{q} = |q|e^{-i\varphi}$, комплексно сопряженное к q , также является корнем характеристического уравнения (6.22), а наряду с комплексными решениями разностного уравнения (6.20)

$$q^n \quad \text{и} \quad \bar{q}^n \quad (6.25)$$

решениями указанного разностного уравнения будут и действительная и мнимая части решений (6.25), т.е.

$$|q|^n \cos n\varphi, \quad |q|^n \sin n\varphi. \quad (6.26)$$

Решения (6.26), как и (6.25), линейно независимы.

Пример 6.1. Найдем общее решение разностного уравнения

$$y_n - 2 \operatorname{ch} \alpha y_{n-1} + y_{n-2} = 0, \quad \alpha \neq 0.$$

Характеристическое уравнение этого разностного уравнения имеет вид

$$q^2 - 2 \operatorname{ch} \alpha q + 1 = 0,$$

а его корни суть

$$q_{1,2} = \operatorname{ch} \alpha \pm \sqrt{\operatorname{ch}^2 \alpha - 1} = e^{\pm \alpha}.$$

Эти корни различные, и поэтому

$$y_n = c_1 e^{\alpha n} + c_2 e^{-\alpha n}.$$

Теорема 6.2. Если q есть корень характеристического уравнения (6.22) кратности $s \geq 1$, то сеточная функция

$$P_{s-1}(n)q^n,$$

где $P_{s-1}(n)$ — произвольный многочлен, степень которого не выше $s-1$, является решением разностного уравнения (6.20). При этом решения

$$n^l q^n, \quad l = 0, \dots, s-1$$

линейно независимы.

Доказательство. Доказательство проведем для случая $s = k = 2$. Покажем сначала, что nq^n есть решение уравнения (6.20) при $k = 2$. Имеем

$$\begin{aligned} & \alpha_0 n q^n + \alpha_1 (n-1) q^{n-1} + \alpha_2 (n-2) q^{n-2} = \\ & = q^{n-2} [(\alpha_0 q^2 + \alpha_1 q + \alpha_2)(n-2) + q(2\alpha_0 q + \alpha_1)] = \\ & = q^{n-1} (2\alpha_0 q + \alpha_1) = 0, \end{aligned}$$

ибо $2\alpha_0 q + \alpha_1 = (\alpha_0 q^2 + \alpha_1 q + \alpha_2)'$ и обязано обращаться в нуль на корне кратности два. Линейная независимость решений q^n и nq^n доказывается так же, как и в предыдущей теореме.

Теорема 6.3. Если правая часть $f(n)$ неоднородного разностного уравнения (6.19) имеет вид $P_m(n)\overset{\circ}{q}^n$, где $P_m(n)$ — многочлен степени m , а $\overset{\circ}{q}$ является s -кратным, $s \geq 0$, корнем характеристического уравнения (6.22), то уравнение (6.19) имеет решение вида

$$y(n) = n^s Q_m(n) \overset{\circ}{q}^n. \quad (6.27)$$

Доказательство смотри, например, в [1].

Упражнение 6.1. Найти общее решение неоднородного разностного уравнения

$$y_n - 2y_{n-1} + y_{n-2} = n(1 + 2^n).$$

Ответ.

$$y_n = c_1 + c_2 n + \frac{1}{6}(n+3)n^2 + (n-2)2^{n+2}.$$

6.4 Системы разностных уравнений

Рассмотрим систему двух линейных однородных разностных уравнений первого порядка с двумя неизвестными функциями

$$\begin{aligned} u_n &= a_{11}u_{n-1} + a_{12}v_{n-1}, \\ v_n &= a_{12}u_{n-1} + a_{22}v_{n-1}, \end{aligned} \quad n \in \mathbb{Z}, \quad (6.28)$$

где a_{ij} , $i, j = 1, 2$ — постоянные. Введем в рассмотрение вектор-функцию

$$y(n) = [u_n \ v_n]^T$$

и матрицу

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix},$$

которую будем предполагать невырожденной, $\det A \neq 0$. Используя введенные обозначения, систему (6.28) можно переписать в векторном виде

$$y(n) = Ay(n-1), \quad \det A \neq 0, \quad \text{или} \quad A^{-1}y(n) = y(n-1). \quad (6.29)$$

В записи (6.29) можно забыть, что $y(n)$ был двумерный вектор, а A — матрица второго порядка. Будем мыслить систему (6.29) как систему m -го порядка. Решение этой системы будем искать в виде

$$y(n) = \xi q^n \quad (6.30)$$

где $q = \text{const} \neq 0$, а ξ — ненулевой m -мерный вектор. Подставляя (6.30) в (6.29), получим

$$\xi q^n = A\xi q^{n-1},$$

а сокращая на $q^{n-1} \neq 0$, приходим к системе

$$\xi q = A\xi.$$

Эта система однородна, и, чтобы у нее были нетривиальные решения, определитель ее матрицы должен быть равен нулю

$$|A - qI| = 0. \quad (6.31)$$

Уравнение (6.31) — характеристическое уравнение системы (6.29) — является алгебраическим уравнением m -ой степени. Если все его корни q_k — собственные значения матрицы A — различны, то соответствующие собственные векторы ξ_k линейно независимы, и общее решение системы (6.29) принимает вид

$$y(n) = \sum_{k=1}^m c_k \xi_k q_k^n. \quad (6.32)$$

Матрица A может иметь полный набор линейно независимых собственных векторов и при наличии кратных корней характеристического уравнения (6.31). И в этом случае общее решение системы (6.29) имеет вид (6.32). Если же у канонической формы матрицы A имеются жордановы клетки, то для отыскания общего решения системы (6.29) нужно поступать так же, как и в случае систем дифференциальных уравнений. На этом мы останавливаться не будем.

Упражнение 6.2. Найти общее решение системы (6.29) с матрицей

$$A = \begin{bmatrix} -1 & -1 \\ 1 & -1 \end{bmatrix}$$

Ответ.

$$y(n) = \left\{ c_1 \begin{bmatrix} -\sin \frac{3\pi n}{4} \\ \cos \frac{3\pi n}{4} \end{bmatrix} + c_2 \begin{bmatrix} \cos \frac{3\pi n}{4} \\ \sin \frac{3\pi n}{4} \end{bmatrix} \right\} q^{n/2}.$$

6.5 Разностная задача на собственные значения

До сих пор мы обсуждали вопросы отыскания общего решения разностного уравнения, которое зависит от k произвольных постоянных, если уравнение имеет порядок k . Чтобы выделить единственное решение разностного уравнения, как и в случае дифференциального уравнения k -го порядка, нужно задать k линейно независимых начальных или граничных условий. Как и для дифференциального уравнения, для разностного уравнения можно поставить задачу на собственные значения.

Займемся этой задачей, решение которой понадобится нам при дальнейших исследованиях. Пусть

$$-y_{n+1} + 2y_n - y_{n-1} = \lambda y_n, \quad n = 1, \dots, N-1, \quad y_0 = y_N = 0. \quad (6.33)$$

Требуется найти такие значения параметра λ (собственные значения), при которых однородная задача (6.33) имеет нетривиальные решения. Если исключить из первого уравнения (6.33) неизвестное y_0 , а из последнего уравнения — неизвестное y_N , то получим обычную алгебраическую задачу на собственные значения для трехдиагональной матрицы

$$A = \begin{bmatrix} 2 & -1 & & & & & \\ -1 & 2 & -1 & & & & \\ & -1 & 2 & -1 & & & \\ \dots & \dots & \dots & \dots & \dots & \dots & \\ & & & -1 & 2 & -1 & \\ & & & & -1 & 2 & \end{bmatrix},$$

порядок которой равен $N-1$. Матрица A симметрична, и потому ее собственные значения действительны, а собственные векторы, отвечающие различным собственным значениям, ортогональны в смысле скалярного произведения

$$(v, w) := \sum_{i=1}^{N-1} v_i w_i.$$

Найдем собственные значения и отвечающие им собственные функции задачи (6.33). Перепишем уравнение (6.33) в виде

$$-y_{n+1} + 2(1 - \lambda/2)y_n - y_{n-1} = 0, \quad n = 1, \dots, N-1$$

и предположим, что

$$|1 - \lambda/2| \leq 1, \quad \text{т.е.} \quad 0 \leq \lambda \leq 4. \quad (6.34)$$

Тогда для некоторого α можно положить

$$1 - \lambda/2 = \cos \alpha/N \quad (6.35)$$

и переписать уравнение так

$$y_{n+1} - 2 \cos \frac{\alpha}{N} y_n + y_{n-1} = 0. \quad (6.36)$$

Характеристическое уравнение, отвечающее уравнению (6.36), есть

$$q^2 - 2 \cos \frac{\alpha}{N} q + 1 = 0,$$

а его корни суть

$$q_{1,2} = \cos \frac{\alpha}{N} \pm \sqrt{\cos^2 \frac{\alpha}{N} - 1} = \cos \frac{\alpha}{N} \pm i \sin \frac{\alpha}{N} = e^{\pm \alpha i/N}.$$

Тем самым,

$$y_n = c_1 \sin \frac{\alpha n}{N} + c_2 \cos \frac{\alpha n}{N} \quad (6.37)$$

есть общее решение уравнения (6.36). Потребуем, чтобы это решение удовлетворяло граничным условиям (6.34). Будем иметь

$$y_0 = c_2 = 0, \quad y_N = c_1 \sin \alpha = 0. \quad (6.38)$$

Поскольку c_1 не может быть нулевым, то отсюда находим, что $\sin \alpha = 0$ и, следовательно,

$$\alpha = k\pi, \quad k \in \mathbb{Z}. \quad (6.39)$$

Из (6.35) находим

$$\lambda = \lambda_k = 2 \left(1 - \cos \frac{k\pi}{N} \right) = 4 \sin^2 \frac{k\pi}{2N}, \quad k \in \mathbb{Z}, \quad (6.40)$$

а из (6.38) —

$$y_n = y_n^{(k)} = c_1 \sin \frac{k\pi n}{N}. \quad (6.41)$$

При $k = 0$ решение $y_n^{(0)} \equiv 0$ и, следовательно, число $\lambda_0 = 0$ не является собственным значением. При $k = N$ решение $y_n^{(N)} = c \sin N\pi \equiv 0$, и λ_N тоже не может быть собственным значением. Собственные значения

$$\lambda_k = 4 \sin^2 \frac{k\pi}{2N}, \quad k = 1, \dots, N-1 \quad (6.42)$$

различны, ибо функция $\sin t$ при $0 < t < \pi/2$ является монотонной. Поскольку изучаемая задача эквивалентна алгебраической задаче на собственные значения для матрицы $(N-1)$ порядка, то соотношения (6.42) задают все собственные значения задачи (6.33). Собственные функции

$$y_n^{(k)} = c_1 \sin \frac{k\pi n}{N}, \quad k = 1, \dots, N-1 \quad (6.43)$$

ортогональны. Подсчитаем их нормы

$$\begin{aligned} \|y_n^{(k)}\|^2 &= (y_n^{(k)}, y_n^{(k)}) = c_1^2 \sum_{n=1}^{N-1} \sin^2 \frac{k\pi n}{N} = \\ &= c_1^2 \sum_{n=1}^{N-1} \frac{1 - \cos \frac{2k\pi n}{N}}{2} = c_1^2 \left[\frac{N-1}{2} - \frac{1}{2} \sum_{n=1}^{N-1} \cos \frac{2k\pi n}{N} \right]. \end{aligned}$$

Далее,

$$\sum_{n=1}^{N-1} \cos \frac{2k\pi n}{N} = \operatorname{Re} \sum_{n=1}^{N-1} \left(e^{\frac{2ik\pi}{N}} \right)^n = \operatorname{Re} \frac{e^{2ik\pi} - e^{2ik\pi/N}}{e^{2ik\pi/N} - 1} = -1.$$

Тем самым,

$$\|y_n^{(k)}\|^2 = c_1^2 N/2 = 1 \quad \text{при} \quad c_1 = \sqrt{2/N}, \quad (6.44)$$

а ортонормированные собственные функции суть

$$y_n^{(k)} = \sqrt{2/N} \sin \frac{k\pi n}{N}, \quad k = 1, \dots, N-1. \quad (6.45)$$

Из (6.42) следует, что для всех собственных значений предположение (6.34) выполнено. Поэтому в рассмотрении противоположного предположения смысла нет.

Упражнение 6.3. Решить следующую задачу на собственные значения

$$\begin{aligned} -y_{n+1} + 2y_n - y_{n-1}, \quad n = 1, \dots, N-1, \\ -y_1 + y_0 + \frac{\lambda}{2}y_0 = 0, \quad y_N - y_{N-1} + \frac{\lambda}{2}y_N = 0. \end{aligned}$$

Ответ.

$$\lambda_k = 4 \sin^2 \frac{k\pi}{2N}, \quad k = 0, \dots, N, \quad y_n^{(k)} = \sqrt{2/N} \cos \frac{k\pi n}{N}.$$

§ 7

Ортогональные многочлены

7.1 Общие ортогональные многочлены

Функцию $\rho(x) \neq 0$ будем называть весовой функцией на интервале $(-1, 1)$, если на этом интервале она неотрицательна и интегрируема.

Пусть на $(-1, 1)$ задана последовательность многочленов

$$P_0(x), P_1(x), \dots, P_n(x), \dots, \quad (7.1)$$

в которой каждый многочлен $P_n(x)$ имеет степень n . Если для любых двух многочленов из этой последовательности выполняется условие

$$(P_m, P_n) := \int_{-1}^1 \rho(x) P_m(x) P_n(x) dx = 0, \quad m \neq n,$$

то многочлены (7.1) называются ортогональными на $(-1, 1)$ с весом $\rho(x)$.

Лемма 7.1. *Если в системе из $(n + 1)$ ортогональных многочленов*

$$P_0(x), P_1(x), \dots, P_n(x)$$

каждый многочлен $P_k(x)$ имеет степень k , то всякий многочлен $Q_n(x)$ степени n можно единственным образом представить в виде

$$Q_n(x) = a_0 P_0(x) + a_1 P_1(x) + \dots + a_n P_n(x). \quad (7.2)$$

Доказательство. Пусть ортогональные многочлены $P_k(x)$ имеют вид

$$P_k(x) = c_0^{(k)} + c_1^{(k)} x + \dots + c_k^{(k)} x^k, \quad c_k^{(k)} \neq 0,$$

а многочлен

$$Q_n(x) = c_0 + c_1 x + \dots + c_n x^n.$$

Подставляя эти представления в (7.2) и приравнивая коэффициенты при одинаковых степенях x^k , получим следующую систему линейных алгебраических уравнений для определения неизвестных коэффициентов a_k

$$\begin{aligned} c_0^{(0)} a_0 + c_0^{(1)} a_1 + c_0^{(2)} a_2 + \dots + c_0^{(n)} a_n &= c_0, \\ c_1^{(1)} a_1 + c_1^{(2)} a_2 + \dots + c_1^{(n)} a_n &= c_1, \\ \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots & \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \\ c_n^{(n)} a_n &= c_n. \end{aligned}$$

По условию теоремы коэффициенты $c_k^{(k)}$ рассматриваемой системы многочленов отличны от нуля, и, следовательно, эта алгебраическая система имеет единственное решение. Лемма доказана.

Замечание 7.1. Умножая соотношение (7.2) на $\rho(x)P_k(x)$ и интегрируя результат по интервалу $(-1, 1)$, легко находим, что

$$a_k = \frac{(P_k, Q_n)}{\|P_k\|^2}, \quad \|P_k\|^2 = (P_k, P_k).$$

Лемма 7.2. Для всякой весовой функции $\rho(x)$ существует единственная последовательность **ортонормированных** многочленов $\{P_n(x)\}$, имеющих положительный коэффициент при старшей степени.

Доказательство. Обозначим коэффициент при старшей степени x многочлена $P_n(x)$ через μ_n . Доказательство теоремы проведем методом полной математической индукции. Имеем, $P_0(x) = \mu_0 > 0$ и, следовательно,

$$(P_0, P_0) = \mu_0^2 (1, 1) = 1.$$

Поэтому

$$\mu_0 = 1/\sqrt{(1, 1)}$$

и многочлен $P_0(x)$ определен.

Пусть определены ортонормированные многочлены

$$P_0(x), P_1(x), \dots, P_{n-1}(x).$$

Определим многочлен $P_n(x)$. Будем его искать в виде $P_n(x) = \mu_n x^n + Q_{n-1}(x)$. В силу леммы 7.1 находим, что

$$P_n(x) = \mu_n x^n + \sum_{k=0}^{n-1} a_k P_k(x),$$

где числа μ_n и a_k подлежат определению. Умножая это соотношение скалярно на $P_m(x)$, $m = 0, \dots, n-1$, находим, что

$$0 = \mu_n (x^n, P_m(x)) + a_m, \quad m = 0, \dots, n-1,$$

т.е.

$$a_m = -\mu_n (x^n, P_m(x))$$

и, следовательно,

$$P_n(x) = \mu_n \left[x^n - \sum_{k=0}^{n-1} (x^n, P_k) x^k \right]$$

есть произвольный ортогональный многочлен степени n . Умножая его скалярно на самого себя и требуя нормированности, находим

$$1 = \mu_n^2 \underbrace{\left(\left(x^n - \sum_{k=1}^{n-1} (x^n, P_k) x^k \right)^2, 1 \right)}_{\underset{\vee}{0}}.$$

Отсюда определяем $\mu_n > 0$. Лемма доказана.

Лемма 7.3. Если $P_n(x)$ принадлежит совокупности ортогональных с весом $\rho(x)$ многочленов, то для всякого многочлена $Q_m(x)$ степени $m < n$

$$(Q_m, P_n(x)) = 0, \quad m < n. \quad (7.3)$$

Доказательство. В силу леммы 7.1

$$Q_m(x) = a_0 P_0(x) + a_1 P_1(x) + \dots + a_m P_m(x).$$

Подставляя это представление в (7.3), получаем утверждение леммы.

Лемма 7.4. Если весовая функция $\rho(x)$ четная, то каждый ортогональный многочлен $P_n(x)$ содержит только те степени x , которые имеют одинаковую с номером n четность, т.е.

$$P_n(-x) \equiv (-1)^n P_n(x). \quad (7.4)$$

Доказательство. Пусть $\rho(x) = \rho(-x)$ и

$$\int_{-1}^1 \rho(x) P_n(x) P_m(x) dx = 0, \quad m = 0, \dots, n-1.$$

Заменой переменной интегрирования $x = -t$ эти условия приводятся к виду

$$\int_{-1}^1 \rho(t) P_n(-t) P_m(-t) dt = 0, \quad m = 0, \dots, n-1,$$

т.е. $P_m(-t)$ тоже ортогональные многочлены. Но в силу леммы 7.2 любой ортогональный многочлен определен с точностью до множителя, и поэтому

$$P_n(-x) = c_n P_n(x).$$

Отсюда в частности следует, что

$$(-1)^n a_n x^n = c_n a_n x^n,$$

т.е. $c_n = (-1)^n$ и поэтому

$$P_n(-x) = (-1)^n P_n(x).$$

Лемма доказана.

Теорема 7.1. *Для любых трех последовательных ортогональных многочленов справедлива рекуррентная формула*

$$\alpha_n P_{n+1}(x) = (x - \beta_n) P_n(x) - \gamma_n P_{n-1}(x). \quad (7.5)$$

Доказательство. Перепишем (7.5) в виде

$$x P_n(x) = \alpha_n P_{n+1}(x) + \beta_n P_n(x) + \gamma_n P_{n-1}(x).$$

В левой части этого равенства стоит многочлен степени $n + 1$. В силу леммы 7.1 он может быть разложен по многочленам P_0, \dots, P_{n+1}

$$x P_n(x) = \sum_{k=0}^{n+1} c_k^{(n+1)} P_k(x), \quad (7.6)$$

где в силу замечания 7.1

$$c_k^{(n+1)} = \frac{(x P_n, P_k)}{\|P_k\|^2}. \quad (7.7)$$

Но тогда с учетом (7.6)

$$\begin{aligned} c_k^{(n+1)} &= \frac{1}{\|P_k\|^2} (P_n, x P_k) = \frac{1}{\|P_k\|^2} \left(P_n, \sum_{j=0}^{k+1} c_j^{(k+1)} P_j \right) = \\ &= \frac{1}{\|P_k\|^2} \sum_{j=0}^{k+1} c_j^{(k+1)} (P_n, P_j). \end{aligned}$$

Отсюда следует, что если $k + 1 < n$, то

$$c_k^{(n+1)} = 0, \quad k + 1 < n$$

и поэтому

$$x P_n(x) = c_{n+1}^{(n+1)} P_{n+1}(x) + c_n^{(n+1)} P_n(x) + c_{n-1}^{(n+1)} P_{n-1}(x), \quad (7.8)$$

т.е.

$$\alpha_n = c_{n+1}^{(n+1)}, \quad \beta_n = c_n^{(n+1)}, \quad \gamma_n = c_{n-1}^{(n+1)}. \quad (7.9)$$

Теорема доказана.

Теорема 7.2. Все нули ортогонального многочлена $P_n(x)$ действительны, различны и расположены на интервале $(-1, 1)$.

Доказательство. Достаточно показать, что многочлен $P_n(x)$ на $(-1, 1)$ меняет знак n раз. Допустим противное, т.е. что многочлен $P_n(x)$ меняет знак только в точках $\xi_1, \xi_2, \dots, \xi_m$, где $m < n$. Тогда многочлен

$$Q_m(x) = (x - \xi_1)(x - \xi_2) \dots (x - \xi_m)$$

также меняет знак точно в этих же самых точках. Следовательно, произведение $P_n(x)Q_m(x) \neq 0$ сохраняет знак на $(-1, 1)$ и следовательно

$$\int_{-1}^1 \rho(x) P_n(x) Q_m(x) dx \neq 0,$$

что противоречит лемме 7.3. Противоречие снимается, если $n = m$. Теорема доказана.

Замечание 7.2. В силу доказанной теоремы для нулей $x_k^{(n)}$ ортогонального многочлена $P_n(x)$ имеют место неравенства

$$-1 < x_1^{(n)} < x_2^{(n)} < \dots < x_k^{(n)} < \dots < x_n^{(n)} < 1. \quad (7.10)$$

7.2 Многочлены Чебышева первого рода

Рассмотрим следующее однородное разностное уравнение второго порядка с постоянными коэффициентами

$$y_n - 2xy_{n-1} + y_{n-2} = 0. \quad (7.11)$$

Здесь x — параметр. Поставим для (7.11) начальные условия

$$y_0 = 1, \quad y_1 = x. \quad (7.12)$$

Тогда

$$\begin{aligned} y_2 &= 2x \cdot x - 1 = 2x^2 - 1, \\ y_3 &= 2x(2x^2 - 1) - x = 4x^3 - 3x, \\ y_4 &= 8x^4 - 8x^2 + 1, \dots \end{aligned} \quad (7.13)$$

Очевидно, что значение решения задачи (7.11), (7.12) в узле n есть многочлен от x степени n .

Найдем решение задачи (7.11), (7.12) в явном виде. Характеристическое уравнение разностного уравнения (7.11) имеет вид

$$q^2 - 2xq + 1 = 0,$$

а его корни суть

$$q_1 = q = x + \sqrt{x^2 - 1} \quad \text{и} \quad q_2 = 1/q. \quad (7.14)$$

Поэтому общее решение уравнения (7.11) есть

$$y_n = c_1 q^n + c_2 q^{-n}.$$

Полагая здесь $n = 0$ и $n = 1$ и принимая во внимание начальные условия (7.12), находим, что

$$\begin{aligned} y_0 &= c_1 + c_2 = 1, \\ y_1 &= c_1 q + c_2 q^{-1} = x. \end{aligned} \quad (7.15)$$

Преобразем второе из уравнений (7.15) с учетом (7.14) и первого уравнения,

$$\begin{aligned} y_1 &= c_1(x + \sqrt{x^2 - 1}) + c_2(x - \sqrt{x^2 - 1}) = \\ &= (c_1 + c_2)x + (c_1 - c_2)\sqrt{x^2 - 1} = x + (c_1 - c_2)\sqrt{x^2 - 1} = x. \end{aligned}$$

Отсюда следует, что $c_1 = c_2$, а с учетом (7.15) находим, что

$$c_1 = c_2 = 1/2,$$

и поэтому

$$y_n = \frac{q^n + q^{-n}}{2} \quad (7.16)$$

есть решение задачи (7.11), (7.12).

Как было замечено раньше, это есть многочлен от x степени n . Пусть $|x| < 1$. Тогда в силу (7.14)

$$q = x + i\sqrt{1 - x^2}$$

и, следовательно, $|q| = 1$. Пусть $q = e^{i\varphi}$. Тогда

$$\begin{aligned} x &= \cos \varphi, \quad \varphi = \arccos x, \\ y_n &= \frac{e^{in\varphi} + e^{-in\varphi}}{2} = \cos n\varphi = \cos[n \arccos x]. \end{aligned} \quad (7.17)$$

Определение 7.1. Алгебраические многочлены

$$T_n(x) = \cos[n \arccos x], \quad |x| < 1, \quad n = 0, 1, \dots \quad (7.18)$$

называются многочленами Чебышева первого рода.

Они принадлежат к семейству ортогональных многочленов. Определим весовую функцию $\rho(x)$, при которой многочлены $T_n(x)$ будут ортогональными. Из (7.17) следует, что

$$d\varphi = -\frac{dx}{\sqrt{1 - x^2}}, \quad x = 1 \text{ при } \varphi = 0 \quad \text{и} \quad x = -1 \text{ при } \varphi = \pi.$$

Принимая теперь во внимание (7.18), находим, что при $m \neq n$

$$0 = \int_0^\pi \cos m\varphi \cos n\varphi d\varphi = \int_{-1}^1 \frac{1}{\sqrt{1-x^2}} T_m(x)T_n(x)dx.$$

Тем самым

$$\rho(x) = (1-x^2)^{-1/2}. \quad (7.19)$$

7.3 Свойства многочленов Чебышева

1°. При четном n многочлен $T_n(x)$ является четной функцией x , а при нечетном n — нечетной.

Доказательство следует из (7.19) и леммы 7.4.

2°. Коэффициент при старшей степени многочлена $T_n(x)$ для $n \geq 1$ равен 2^{n-1} , т.е. $\mu_n = 2^{n-1}$, а

$$T_n(x) = 2^{n-1}x^n + \dots$$

Доказательство. См. рекуррентную формулу (7.11).

3°. Нули многочлена $T_n(x)$ расположены в точках

$$x_k = -\cos \frac{(2k-1)\pi}{2n}, \quad k = 1, 2, \dots, n. \quad (7.20)$$

Доказательство. Из (7.18) находим, что

$$n \arccos x_k = -\frac{\pi}{2} + k\pi = \frac{(2k-1)\pi}{2}$$

или

$$\arccos x_k = \frac{(2k-1)\pi}{2n},$$

т.е.

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

Так как функции $T_n(x)$ являются либо четными, либо нечетными, то нули $T_n(x)$ расположены симметрично относительно начала координат

$$x_{n+1-k} = -x_k = -\cos \frac{(2k-1)\pi}{2n}.$$

Перенумеровывая нули в обратном порядке, приходим к (7.20).

4°. $\max_{[-1,1]} |T_n(x)| = 1$, причем

$$T_n(x_m) = (-1)^m,$$

где

$$x_m = -\cos \frac{m\pi}{n}, \quad m = 0, \dots, n. \quad (7.21)$$

Доказательство очевидно.

5°. Среди всех многочленов степени n с единичным коэффициентом при старшей степени многочлен

$$\bar{T}_n(x) = \frac{1}{2^{n-1}} T_n(x), \quad n \geq 1$$

на $[-1, 1]$ имеет наименьшее значение максимума модуля.

Доказательство. Допустим противное, т.е. допустим существование такого многочлена $\bar{P}_n(x) = x^n + \dots$, что

$$\max_{[-1,1]} |\bar{P}_n(x)| < \max_{[-1,1]} |\bar{T}_n(x)|. \quad (7.22)$$

Тогда $\bar{T}_n(x) - \bar{P}_n(x) \neq 0$ и это есть многочлен степени не выше $(n-1)$. Более того, в $(n+1)$ точке (7.21) этот многочлен принимает отличные от нуля значения с чередующимися знаками.

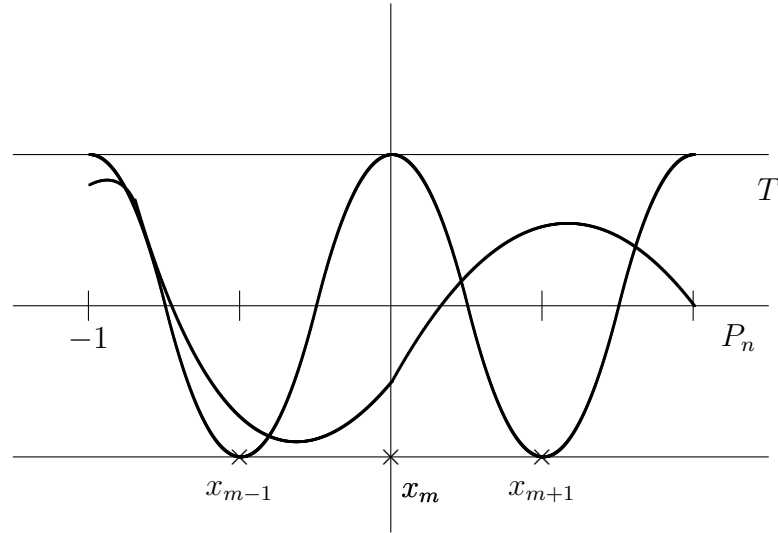


Рис. 1

Но это означает, что алгебраический многочлен $\bar{T}_n(x) - \bar{P}_n(x)$ степени меньшей n обращается в нуль по крайней мере в n точках, что невозможно.

Замечание 7.3. Можно доказать, что если $\bar{P}_n(x) = x^n + \dots$, $n \geq 1$, и

$$\max_{[-1,1]} |\bar{P}_n(x)| = 2^{-n+1},$$

то $\bar{P}_n(x) \equiv \bar{T}_n(x) = 2^{-n-1} T_n(x)$.

Благодаря свойству 5° многочлены Чебышева $T_n(x)$ называются многочленами, наименее уклоняющимися от нуля.

6°. Если $x > 1$, то

$$T_n(x) = \operatorname{ch} n \operatorname{Arch} x,$$

где

$$\operatorname{Arch} x = \ln(x + \sqrt{x^2 - 1}).$$

Доказательство. В силу (7.16)

$$T_n(x) = \frac{q^n + q^{-n}}{2} = \frac{e^{n \ln q} + e^{-n \ln q}}{2} = \operatorname{ch} n \ln q = \operatorname{ch} n \ln(x + \sqrt{x^2 - 1}).$$

Замечание 7.4. $\operatorname{Arch} x$ — обратная функция к $\operatorname{ch} x$.

Упражнение 7.1. Доказать, что

$$\operatorname{ch} \operatorname{Arch} x = x, \quad \operatorname{Arch} \operatorname{ch} x = x.$$

7.4 Многочлены Лежандра

Многочленами Лежандра называются многочлены, которые ортогональны друг другу на $[-1, 1]$ с весом $\rho \equiv 1$. Обозначаются они через $P_n(x)$

$$\int_{-1}^1 P_m(x) P_n(x) dx = 0, \quad m \neq n.$$

Если $P_0(x) \equiv 1$, то $P_1(x) \equiv x$.

Трехточечное рекуррентное соотношение для многочленов Лежандра имеет вид

$$(n+1)P_{n+1}(x) - (2n+1)xP_n(x) + nP_{n-1}(x) = 0$$

и следовательно

$$P_2(x) = \frac{1}{2}(3x^2 - 1), \quad P_3(x) = \frac{1}{2}(5x^3 - 3x),$$

$$P_4(x) = \frac{1}{8}(35x^4 - 30x^2 + 3),$$

$$P_5(x) = \frac{1}{8}(63x^5 - 70x^3 + 15x)$$

и т.д.

§ 8

Итерационные методы

В предыдущих лекциях для системы линейных алгебраических уравнений

$$Ax = b \tag{8.1}$$

с квадратной невырожденной матрицей были рассмотрены четыре прямых метода отыскания решения:

- а) метод Гаусса (LU -разложение, треугольное разложение) и его модификация с выбором ведущего элемента,
- б) метод Холецкого, применяемый в случае симметричной положительно определенной матрицы,
- в) метод вращений,
- г) метод отражений.

Все эти методы позволяют в принципе (при отсутствии ошибок округления) найти точное решение за конечное число действий. Это число действий было оценено нами величиной $O(n^3)$, где n — порядок системы. Если матрица A системы имеет ленточную структуру с полушириной ленты p много меньшей n , то ленточные варианты первых двух методов позволяют найти точное решение с меньшей, чем $O(n^3)$, затратой действий.

В этой лекции мы рассмотрим другой класс методов решения системы (8.1) — итерационных. Эти методы, как правило, если и позволяют найти точное решение системы (8.1), то только как предел при стремлении числа итераций (а, следовательно, и действий) к бесконечности. Однако для широкого класса задач, встречающихся в приложениях, те или иные итерационные методы могут оказаться предпочтительнее с точки зрения используемых трудозатрат, чем описанные прямые.

8.1 Одношаговые итерационные методы

Из курса "Введение в численные методы" известно, что многие одношаговые итерационные методы могут быть записаны в так называемой канонической форме

$$B \frac{x^{k+1} - x^k}{\tau} + Ax^k = b, \quad k = 0, 1, \dots, \quad (8.2)$$

где B — некоторая матрица, определяющая итерационный метод, а τ — итерационный параметр. В частности, в виде (8.2) могут быть записаны метод Якоби, метод Гаусса-Зейделя, метод последовательной верхней релаксации, метод простых итераций. Предположим, что

$$A = A^T > 0, \quad B = B^T > 0. \quad (8.3)$$

При этих предположениях в курсе "Введение в численные методы" доказано, что если

$$B > \frac{\tau}{2}A, \quad (8.4)$$

т.е. если для любого ненулевого вектора x справедливо $(Bx, x) > \frac{\tau}{2}(Ax, x)$, то итерационный метод (8.2) является сходящимся.

Для метода простых итераций, который имеет вид (8.2) с $B = I$, кроме того, дана и оценка скорости сходимости. Именно, если

$$\tau = 2/(\lambda_1 + \lambda_n), \quad (8.5)$$

где λ_1 и λ_n , соответственно, наименьшее и наибольшее собственные значения матрицы A , то для погрешности итераций справедлива оценка

$$\|x - x^k\| \leq q \|x - x^{k-1}\|, \quad q = \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} = \frac{1 - \lambda_1/\lambda_n}{1 + \lambda_1/\lambda_n} < 1, \quad (8.6)$$

где $\|\cdot\|$ — евклидова длина.

Из (8.6) следует, что если $\lambda_n \gg \lambda_1$, то число q очень близко к единице, а скорость сходимости итераций очень низкая. Но при указанном выборе нормы вектора

$$\|A\| = \max_i \lambda_i(A) = \lambda_n, \quad \|A^{-1}\| = \max_i \lambda_i(A^{-1}) = \lambda_1^{-1}$$

и

$$\lambda_n/\lambda_1 = \|A\| \|A^{-1}\| = \text{cond } A := \varkappa(A) := \varkappa.$$

Таким образом, если матрица A плохо обусловлена, а это типичная ситуация, то метод простых итераций будет сходиться очень медленно.

8.2 Неявные методы

Какие есть пути увеличения скорости сходимости итерационных методов? Изучим влияние матрицы B из (8.2) на скорость сходимости. В силу (8.3) существует матрица $B^{1/2}$ такая, что

$$B^{1/2} = (B^{1/2})^T > 0 \quad \text{и} \quad B^{1/2} B^{1/2} = B.$$

Эта матрица называется квадратным корнем из матрицы B .

Напомним построение матрицы $B^{1/2}$. Пусть λ — собственное значение матрицы B , а ξ — отвечающий ему собственный вектор, т.е. $B\xi = \lambda\xi$. Перенумеруем все собственные значения матрицы B и введем в рассмотрение диагональную матрицу $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, образованную этими собственными значениями, и ортогональную матрицу $\Xi = [\xi_1 \xi_2 \dots \xi_n]$, образованную ортонормированными собственными векторами ξ_i матрицы B , упорядоченными в соответствии с нумерацией собственных значений. Поскольку $\Xi\Lambda = [\lambda_1 \xi_1 \lambda_2 \xi_2 \dots \lambda_n \xi_n]$, то $B\Xi = \Xi\Lambda$. Отсюда следует, что $B = \Xi\Lambda\Xi^T$.

Очевидно, что матрица $\Lambda^{1/2} = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_n})$. Поэтому $B^{1/2} = \Xi\Lambda^{1/2}\Xi^T$.

Введем обозначение

$$B^{1/2}x^k = y^k, \quad x^k = (B^{1/2})^{-1}y^k = B^{-1/2}y^k. \quad (8.7)$$

Тогда (8.2) можно переписать так

$$B^{1/2} \frac{y^{k+1} - y^k}{\tau} + AB^{-1/2}y^k = b,$$

а после применения к этому соотношению матрицы $B^{-1/2}$, получим

$$\frac{y^{k+1} - y^k}{\tau} + B^{-1/2}AB^{-1/2}y^k = B^{-1/2}b =: f. \quad (8.8)$$

Обозначая

$$B^{-1/2}AB^{-1/2} = C, \quad (8.9)$$

будем иметь соотношения

$$\frac{y^{k+1} - y^k}{\tau} + Cy^k = f, \quad k = 0, 1, \dots, \quad (8.10)$$

которые по форме полностью совпадают с (8.2). Очевидно, что $C = C^T > 0$, и поэтому для итерационного метода (8.10) (а это есть метод простых итераций) при

$$\tau = \frac{2}{\lambda_1(C) + \lambda_n(C)} \quad (8.11)$$

справедлива оценка

$$\|y^{k+1} - y\| \leq q(C) \|y^k - y\|, \quad q(C) = \frac{\lambda_n(C) - \lambda_1(C)}{\lambda_n(C) + \lambda_1(C)} < 1, \quad (8.12)$$

где $\lambda_1(C)$ и $\lambda_n(C)$ — соответственно минимальное и максимальное собственные значения матрицы C :

$$C\xi = \lambda(C)\xi. \quad (8.13)$$

Здесь ξ — собственные векторы матрицы C . Подставляя в (8.13) представление C из (8.9), получим

$$B^{-1/2}AB^{-1/2}\xi = \lambda(C)\xi,$$

а, обозначая

$$B^{-1/2}\xi = \eta, \quad \xi = B^{1/2}\eta$$

и применяя к последней задаче матрицу $B^{1/2}$, будем иметь

$$A\eta = \lambda(C)B\eta. \quad (8.14)$$

Таким образом, $\lambda_1(C)$ и $\lambda_n(C)$ суть минимальное и максимальное собственные значения обобщенной задачи на собственные значения (8.14).

Подставляя в (8.12) y^k из (8.7) и обозначая $(Bx, x) = \|x\|_B^2$, получим

$$\|x^k - x\|_B \leq q(C) \|x^{k-1} - x\|_B.$$

Итак, если $B = B^T > 0$ такова, что

$$\lambda_n/\lambda_1 > \lambda_n(C)/\lambda_1(C),$$

то итерационный метод (8.2), (8.11) с этой матрицей B будет сходиться быстрее (в смысле B -нормы), чем метод простых итераций. Наибольшую скорость сходимости мы получим, выбирая $B = A$. Тогда

$$C = I, \quad \lambda_1(C) = \lambda_n(C) = 1, \quad \tau = 1$$

и (8.2) принимает вид

$$Ax^k = b,$$

что с точностью до обозначений совпадает с (8.1). Метод сходится за одну итерацию. Лучшего быть не может. Но мы пришли к тому, от чего хотели уйти: нам снова нужно решать систему (8.1). Отсюда следует, что на выбор матрицы B нужно наложить весьма серьезные ограничения — матрица B должна быть относительно легко обратима. Таковыми, например, являются диагональные и треугольные матрицы. Хотя последние и не являются симметричными, это неудобство легко устранить, выбирая в качестве B подходящее произведение треугольных матриц. Однако скорость сходимости метода (8.2) при таком выборе B из очевидных соображений остается сравнительно низкой.

В настоящее время неизвестно регулярных способов хорошего выбора матрицы B для произвольной A . Все удачные находки так или иначе связаны со спецификой матрицы A .

8.3 Нестационарные итерационные методы

Другой путь увеличения скорости сходимости итерационного метода (8.2) состоит в том, чтобы вместо одного итерационного параметра τ использовать несколько — свой итерационный параметр на каждой итерации. Итерационные методы такого типа называются нестационарными и имеют вид

$$B \frac{x^{k+1} - x^k}{\tau_{k+1}} + Ax^k = b, \quad k = 0, 1, \dots \quad (8.15)$$

или с учетом (8.7), (8.10)

$$\frac{y^{k+1} - y^k}{\tau_{k+1}} + Cy^k = f, \quad k = 0, 1, \dots \quad (8.16)$$

Укажем один из возможных способов выбора итерационных параметров τ_k . Введем обозначение

$$z^k = y^k - y,$$

где

$$Cy = f, \quad (8.17)$$

и вычтем (8.17) из (8.16). В результате получим задачу для z^k :

$$\frac{z^{k+1} - z^k}{\tau_{k+1}} + Cz^k = 0, \quad k = 0, 1, \dots, \quad z^0 = y^0 - y. \quad (8.18)$$

Отсюда

$$z^{k+1} = (I - \tau_{k+1}C)z^k,$$

и, следовательно,

$$z^k = \prod_{j=1}^k (I - \tau_j C) z^0, \quad (8.19)$$

т.е.

$$z^k = P_k(C)z^0,$$

где

$$P_k(t) = \prod_{j=1}^k (1 - \tau_j t) = 1 + a_1^{(k)}t + \dots + a_k^{(k)}t^k. \quad (8.20)$$

Из (8.19) находим, что

$$\|z^k\| = \|y^k - y\| = \left\| \prod_{j=1}^k (I - \tau_j C) z^0 \right\| \leq \left\| \prod_{j=1}^k (I - \tau_j C) \right\| \|z^0\|. \quad (8.21)$$

Но

$$\left\| \prod_{j=1}^k (I - \tau_j C) \right\| = \max_l |\lambda_l(P_k(C))| = \max_l |P_k(\lambda_l(C))|. \quad (8.22)$$

Поскольку мы хотим, чтобы итерации сходились как можно быстрее, то можно поставить задачу о минимизации $\|P_k(C)\|$ в зависимости от итерационных параметров τ_j , $j = \overline{1, k}$. В силу (8.20), (8.22) эта задача эквивалентна задаче построения многочлена $P_k(t)$ степени k с единичным свободным членом, который в точках спектра матрицы C наиболее близок к нулю. Но поставленная задача практически не разрешима. Однако вместо нее можно поставить близкую задачу о построении $P_k(t)$, наименее отклоняющегося от нуля не на спектре, а на отрезке $[\lambda_1, \lambda_n]$, где этот спектр расположен. Эта задача много проще, и решение ее известно. Найдем это решение.

Итак, среди многочленов степени k таких, что $Q_k(0) = 1$ (см. (8.20)), требуется найти многочлен $P_k(t)$, максимум модуля которого на $[\lambda_1, \lambda_n]$ минимален.

Линейной заменой переменной $t = ax + b$ переведем отрезок $[\lambda_1, \lambda_n]$ в отрезок $[1, -1]$. Имеем

$$\left. \begin{aligned} \lambda_1 &= a + b \\ \lambda_n &= -a + b \end{aligned} \right\}, \quad b = \frac{\lambda_1 + \lambda_n}{2}, \quad a = \frac{\lambda_1 - \lambda_n}{2},$$

т.е.

$$t = \frac{\lambda_n + \lambda_1}{2} - \frac{\lambda_n - \lambda_1}{2} x = \frac{\lambda_1 + \lambda_n}{2} \left[1 - \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} x \right] = \frac{1}{\tau_0} [1 - \rho_0 x], \quad (8.23)$$

где

$$\tau_0 = \frac{2}{\lambda_1 + \lambda_n}, \quad \rho_0 = \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} < 1. \quad (8.24)$$

Заметим, что τ_0 совпадает с τ из (8.5), (8.11) для стационарного итерационного процесса, а ρ_0 совпадает с q из (8.6), (8.12) и характеризует скорость сходимости этого процесса.

Пусть

$$P_k(t) = \widehat{P}_k(x). \quad (8.25)$$

Поскольку $t = 0$ отвечает точка $x = \rho_0^{-1}$ (см. (8.23)), то должно быть

$$P_k(0) = \widehat{P}_k\left(\frac{1}{\rho_0}\right) = 1. \quad (8.26)$$

Наша задача свелась к отысканию многочлена $\widehat{P}_k(x)$, наименее отклоняющегося от нуля на $[-1, 1]$ и удовлетворяющего условию (8.26). С похожей задачей мы уже сталкивались на прошлой лекции. Мы знаем, что среди многочленов степени k вида $x^k + \dots$ наименее отклоняется от нуля на $[-1, 1]$ многочлен

$$\overline{T}_k(x) = \frac{1}{2^{k-1}} T_k(x),$$

где $T_k(x)$ — многочлен Чебышева первого рода. Разумеется, сам многочлен Чебышева $T_k(x)$ является наименее отклоняющимся от нуля на $[-1, 1]$ среди многочленов вида $P_k(x) = 2^{k-1}x^k + \dots$.

Пусть $T_k\left(\frac{1}{\rho_0}\right) = \frac{1}{q_k}$. Тогда очевидно, что искомое решение дает многочлен

$$\widehat{P}_k(x) = q_k T_k(x). \quad (8.27)$$

Из свойства 6° многочленов T_k

$$q_k = \left[\operatorname{ch} k \operatorname{Arch} \frac{1}{\rho_0} \right]^{-1} < 1. \quad (8.28)$$

Получим еще одно представление для q_k . Из (8.28)

$$\operatorname{ch} k \operatorname{Arch} \frac{1}{\rho_0} = \frac{1}{q_k}.$$

Отсюда

$$k \operatorname{Arch} \frac{1}{\rho_0} \equiv k \ln \frac{1 + \sqrt{1 - \rho_0^2}}{\rho_0} = \operatorname{Arch} \frac{1}{q_k} \equiv \ln \frac{1 + \sqrt{1 - q_k^2}}{q_k}.$$

Пусть

$$\begin{aligned} \rho_1 &= \frac{\rho_0}{1 + \sqrt{1 - \rho_0^2}} = \frac{\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}}{1 + \sqrt{1 - \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^2}} = \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1 + 2\sqrt{\lambda_1 \lambda_n}} = \\ &= \frac{\sqrt{\lambda_n} - \sqrt{\lambda_1}}{\sqrt{\lambda_n} + \sqrt{\lambda_1}} = \frac{1 - \sqrt{\lambda_1/\lambda_n}}{1 + \sqrt{\lambda_1/\lambda_n}}. \end{aligned} \quad (8.29)$$

Тогда

$$\ln \frac{1 + \sqrt{1 - q_k^2}}{q_k} = k \ln \frac{1}{\rho_1} = \ln \rho_1^{-k}$$

и

$$\frac{1 + \sqrt{1 - q_k^2}}{q_k} = \rho_1^{-k}.$$

Отсюда, переходя к квадратному уравнению относительно q_k^2 , находим, что

$$q_k = \frac{2\rho_1^{-k}}{1 + \rho_1^{-2k}} = \frac{2\rho_1^k}{1 + \rho_1^{2k}}. \quad (8.30)$$

Итак, мы нашли такой многочлен

$$\widehat{P}_k(x) = q_k T_k(x),$$

что

$$\max_{[-1,1]} |\widehat{P}_k(x)| = q_k$$

и, следовательно, в силу (8.21), (8.7) имеем оценку сходимости

$$\|z^k\| = \|x - x^k\|_B \leq q_k \|x - x^0\|_B. \quad (8.31)$$

Получим формулы для итерационных параметров τ_j . Из (8.25), (8.27) и (8.23) следует, что нули полиномов $P_k(t)$ и $T_k\left(\frac{1-\tau_0 t}{\rho_0}\right)$ совпадают. Так как полином $P_k(t)$ имеет нули в точках $t = 1/\tau_j$, $j = \overline{1, k}$, а нулями полинома Чебышева $T_k(x)$ являются числа (7.21)

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

то с учетом (8.23) находим, что

$$\tau_j = \frac{\tau_0}{1 + \rho_0 \mu_j}, \quad j = \overline{1, k}, \quad (8.32)$$

где

$$\mu_j \in \mathfrak{M}_k = \left\{ -\cos \frac{2i-1}{2k} \pi, \quad i = \overline{1, k} \right\}. \quad (8.33)$$

Из полученной формулы для итерационных параметров τ_j видно, что для их вычисления требуется задать число итераций k . Как это сделать? Обычно в качестве условия окончания итерационного процесса берется неравенство

$$\|z^k\|_B \leq \varepsilon \|z^0\|_B$$

и числом итераций называется наименьшее из чисел k , для которых это неравенство выполняется. Из (8.31) следует, что для рассматриваемого итерационного метода число итераций находится из неравенства $q_k \leq \varepsilon$. Решим это неравенство, используя (8.30):

$$\frac{2\rho_1^k}{1 + \rho_1^{2k}} \leq \varepsilon.$$

Это неравенство эквивалентно следующему неравенству

$$\left(\rho_1^k - \frac{1 + \sqrt{1 - \varepsilon^2}}{\varepsilon} \right) \left(\rho_1^k - \frac{1 - \sqrt{1 - \varepsilon^2}}{\varepsilon} \right) \geq 0.$$

Поскольку $\rho_1 < 1$, то первый сомножитель отрицателен, и должно выполняться условие

$$\rho_1^k \leq \frac{1 - \sqrt{1 - \varepsilon^2}}{\varepsilon} = \frac{\varepsilon}{1 + \sqrt{1 - \varepsilon^2}}.$$

Отсюда

$$k \geq \frac{\ln \frac{1 + \sqrt{1 - \varepsilon^2}}{\varepsilon}}{\ln 1/\rho_1}.$$

Поскольку ε обычно мало, то пользуются следующей формулой

$$k \geq \frac{\ln 2/\varepsilon}{\ln 1/\rho_1}.$$

Сравним скорости сходимости построенного нестационарного итерационного метода с оптимальным методом простых итераций. Из (8.6) следует, что для оптимального метода простых итераций

$$\|x - x^k\| \leq q \|x - x^{k-1}\|,$$

где

$$q = \frac{1 - \lambda_1/\lambda_n}{1 + \lambda_1/\lambda_n} = \rho_0$$

и, следовательно,

$$\|x - x^k\| \leq q^k \|x - x^0\|.$$

Если мы и здесь потребуем, чтобы

$$\|x - x^k\| \leq \varepsilon \|x - x^0\|,$$

то для числа итераций k будем иметь

$$k \ln \frac{1}{\rho_0} \geq \ln \frac{1}{\varepsilon},$$

т.е.

$$k \geq \frac{\ln 1/\varepsilon}{\ln 1/\rho_0}. \quad (8.34)$$

Для плохо обусловленной матрицы A

$$\lambda_1/\lambda_n = [\text{cond } A]^{-1} \ll 1,$$

и поэтому

$$\frac{1}{\rho_0} = 1 + 2(\text{cond } A)^{-1} + O((\text{cond } A)^{-2}).$$

Отсюда и из (8.34)

$$k \approx \frac{1}{2} \text{cond } A \ln 1/\varepsilon. \quad (8.35)$$

Для оптимального итерационного метода будем иметь

$$\rho_1^{-1} = \frac{1 + \sqrt{(\text{cond } A)^{-1}}}{1 - \sqrt{(\text{cond } A)^{-1}}} = 1 + 2\sqrt{(\text{cond } A)^{-1}} + O((\text{cond } A)^{-1})$$

и

$$k \approx \frac{1}{2} (\text{cond } A)^{1/2} \ln 2/\varepsilon.,$$

что много лучше, чем (8.35).

8.4 Об устойчивости

К сожалению, вычисления по формулам (8.15) при произвольном использовании итерационных параметров не являются устойчивыми.

$M_0 = 10^{-p}$ — машинный нуль.

$M_\infty = 10^p$ — бесконечность.

$$10^{3p/4}, \quad 10^{p/2}, \quad 10^{p/4}, \quad 10^{-p/2}, \quad 10^{-3p/4}$$

$$\Pi = 10^{p/4} \in [M_0, M_\infty]$$

$$\underbrace{10^{3p/4} \cdot 10^{p/2}} \cdot 10^{p/4} \cdot \underbrace{10^{-p/2} \cdot 10^{-3p/4}} = M_\infty$$

$$\underbrace{10^{-p/2} \cdot 10^{-3p/4}} \cdot 10^{p/4} \cdot \underbrace{10^{3p/4} \cdot 10^{p/2}} = M_0$$

$$10^{-3p/4} \ 10^{p/2} \ 10^{3p/4} \ 10^{-p/2} \ 10^{p/4}$$

§ 9

Итерационные методы вариационного типа

9.1 Метод скорейшего спуска

Вновь обратимся к решению системы

$$Ax = b, \quad A = A^T > 0. \quad (9.1)$$

Будем для простоты использовать явный нестационарный метод

$$\frac{x^{k+1} - x^k}{\tau_{k+1}} + Ax^k = b, \quad k = 0, 1, \dots \quad (9.2)$$

В предыдущей лекции был указан способ выбора итерационных параметров τ_k , использующий априорную информацию о расположении спектра матрицы A . Сейчас мы рассмотрим другой способ выбора этих параметров.

Пусть, как обычно, $z^k = x^k - x$. Тогда

$$z^{k+1} = z^k - \tau_{k+1}Az^k. \quad (9.3)$$

Вычислим A -норму погрешности z^{k+1} и выразим ее через z^k . Используя (9.3), находим, что

$$\begin{aligned} \|z^{k+1}\|_A^2 &:= (Az^{k+1}, z^{k+1}) = (A(z^k - \tau_{k+1}Az^k), z^k - \tau_{k+1}Az^k) = \\ &= \|z^k\|_A^2 - 2\tau_{k+1}(Az^k, Az^k) + \tau_{k+1}^2(A^2z^k, Az^k). \end{aligned}$$

Выберем τ_{k+1} из условия минимума $\|z^{k+1}\|_A^2$. Дифференцируя по τ_{k+1} и приравнявая производную нулю, найдем, что

$$-2(Az^k, Az^k) + 2\tau_{k+1}(A^2z^k, Az^k) = 0,$$

т.е.

$$\tau_{k+1} = \frac{(Az^k, Az^k)}{(A^2z^k, Az^k)}. \quad (9.4)$$

Казалось бы, что это соотношение не позволяет найти интересующий нас параметр τ_{k+1} , поскольку $z^{k+1} = x^k - x$ не известна. Но нам она и не нужна. Нам нужна

$$Az^k = Ax^k - Ax = Ax^k - b = -r^k.$$

Величина r^k называется невязкой. Тем самым,

$$\tau_{k+1} = \frac{(r^k, r^k)}{(Ar^k, r^k)}, \quad r^k = b - Ax^k. \quad (9.5)$$

Метод (9.2), (9.5) называется методом скорейшего спуска.

Дадим геометрическую интерпретацию этого метода, которая и объяснит его название. Пусть

$$J(x) = \frac{1}{2}(Ax, x) - (b, x) \quad (9.6)$$

— квадратичная функция n переменных x_1, x_2, \dots, x_n . Поставим задачу об отыскании точки минимума этой функции. Для решения этой задачи нужно найти первые производные (9.6) по x_1, x_2, \dots, x_n и приравнять их нулю. Это и будут уравнения для нахождения точки минимума. Перепишем (9.6) в координатном виде

$$J(x) = \frac{1}{2} \sum_{k,j=1}^n a_{kj}x_kx_j - \sum_{k=1}^n b_kx_k$$

и продифференцируем по x_i

$$\frac{\partial J(x)}{\partial x_i} = \frac{1}{2} \sum_{j=1}^n a_{ij}x_j + \frac{1}{2} \sum_{k=1}^n a_{ki}x_k - b_i = \sum_{j=1}^n a_{ij}x_j - b_i, \quad i = 1, \dots, n.$$

Замечание 9.1. $\text{grad } J = Ax - b$.

Из математического анализа известно, что функция наиболее быстро убывает в направлении антиградиента.

Итак, задача отыскания точки минимума функции (9.6) эквивалентна решению системы (9.1) с симметричной матрицей. Если мы найдем способ приближенного нахождения точки минимума функции (9.6), то мы будем иметь метод приближенного нахождения решения системы (9.1).

Построим метод минимизации (9.6). Применительно к (9.6)

$$\nabla J(x) = Ax - b,$$

и процесс минимизации принимает вид

$$x^{k+1} = x^k - \alpha \nabla J(x^k) = x^k - \alpha(Ax^k - b) = x^k + \alpha r^k \quad (9.7)$$

или

$$\frac{x^{k+1} - x^k}{\alpha} + Ax^k = b,$$

что совпадает с (9.2) с точностью до обозначения итерационного параметра. Для определения значения α рассмотрим

$$J(x^{k+1}) = J(x^k - \alpha \nabla J(x^k)) \quad (9.8)$$

как функцию α и найдем такое значение α , при котором J принимает наименьшее значение. В нашем случае

$$\begin{aligned} J(x^k - \alpha \nabla J(x^k)) &= \frac{1}{2} (A(x^k - \alpha(Ax^k - b)), x^k - \alpha(Ax^k - b)) - (b, x^k - \alpha(Ax^k - b)) = \\ &= \frac{1}{2} (A(x^k + \alpha r^k), x^k + \alpha r^k) - (b, x^k + \alpha r^k). \end{aligned}$$

Дифференцируя по α и приравнивая производную нулю, находим, что

$$(A(x^k + \alpha r^k), r^k) - (b, r^k) = 0,$$

откуда

$$\alpha = \alpha_{k+1} = \frac{(r^k, r^k)}{(Ar^k, r^k)},$$

что совпадает с ранее полученным значением (9.5) итерационного параметра. Тем самым, оба метода совпадают, а второй метод дает им название.

Имеет место

Теорема 9.1. *Итерации по методу скорейшего спуска (9.2), (9.5) сходятся не медленнее, чем в оптимальном методе простых итераций. Именно*

$$\|x^k - x\|_A \leq \left(\frac{1 - \lambda_1/\lambda_n}{1 + \lambda_1/\lambda_n} \right)^k \|x^0 - x\|_A,$$

где λ_1 и λ_n суть минимальное и максимальное собственные значения матрицы A .

Доказательство. Из (9.3)

$$\|z^{k+1}\|_A = \|(I - \tau_{k+1}A)z^k\|_A,$$

причем правая часть принимает минимальное значение именно при τ_{k+1} из (9.5). Тем самым, при любом другом значении τ_{k+1} правая часть будет только больше, и, следовательно,

$$\|z^{k+1}\|_A^2 \leq \|(I - \tau A)z^k\|_A^2$$

для любого τ и, в частности, для

$$\tau = \frac{2}{\lambda_1 + \lambda_n}$$

— итерационного параметра метода простых итераций.

Но

$$\begin{aligned} \|(I - \tau A)z^k\|_A^2 &= ((I - \tau A)z^k, (I - \tau A)Az^k) = \\ &= ((I - \tau A)A^{1/2}z^k, (I - \tau A)A^{1/2}z^k) = \\ &= \|(I - \tau A)A^{1/2}z^k\|^2 \leq \|I - \tau A\| \|z^k\|_A^2, \end{aligned}$$

а в силу (8.22), (8.27), (8.28), при $k = 1$

$$\|I - \tau A\| \leq \max_{\lambda \in [\lambda_1, \lambda_n]} |1 - \tau\lambda| = q_1 = \rho_0 = \frac{1 - \lambda_1/\lambda_n}{1 + \lambda_1/\lambda_n}.$$

Собирая оценки для всех k , получим утверждение теоремы.

Из теоремы 9.1 следует, что нестационарный метод (9.2), (9.5) сравним по скорости сходимости с методом простых итераций, и, казалось бы, мы с этим методом не продвинулись вперед. Однако, у этих методов имеется существенное различие. Для использования метода простых итераций требуется информация о границах спектра матрицы A . В случае же метода (9.2), (9.5) такая информация не требуется.

9.2 Неулучшаемость оценки

Покажем на примере, что полученная оценка сходимости достигается, если начальное приближение задано специальным образом. Допустим, что x^0 таково, что

$$x^0 - x = z^0 = c \left(\sqrt{\frac{\lambda_n}{\lambda_1}} \xi_1 + \sqrt{\frac{\lambda_1}{\lambda_n}} \xi_n \right),$$

где λ_1 и λ_n суть минимальное и максимальное собственные значения матрицы A из (9.1), а ξ_1 и ξ_n — отвечающие им ортонормированные собственные векторы. Тогда

$$\begin{aligned} Az^0 &= c \sqrt{\lambda_1 \lambda_n} (\xi_1 + \xi_n), \\ A^2 z^0 &= c \sqrt{\lambda_1 \lambda_n} (\lambda_1 \xi_1 + \lambda_n \xi_n), \\ (Az^0, Az^0) &= c^2 \lambda_1 \lambda_n 2, \\ (A^2 z^0, A^2 z^0) &= c^2 \lambda_1 \lambda_n (\lambda_1 + \lambda_n). \end{aligned}$$

В силу (9.4)

$$\tau_1 = \frac{(Az^0, Az^0)}{(A^2 z^0, A^2 z^0)} = \frac{2}{\lambda_1 + \lambda_n},$$

а в силу (9.3)

$$\begin{aligned}
 z^1 &= c \sqrt{\frac{\lambda_n}{\lambda_1}} \xi_1 + c \sqrt{\frac{\lambda_1}{\lambda_n}} \xi_n - \frac{2}{\lambda_1 + \lambda_n} c \sqrt{\lambda_1 \lambda_n} (\xi_1 + \xi_n) = \\
 &= c \left[\sqrt{\frac{\lambda_n}{\lambda_1}} \xi_1 \left(1 - \frac{2\lambda_1}{\lambda_1 + \lambda_n} \right) + \sqrt{\frac{\lambda_1}{\lambda_n}} \xi_n \left(1 - \frac{2\lambda_n}{\lambda_1 + \lambda_n} \right) \right] = \\
 &= c \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \left[\sqrt{\frac{\lambda_n}{\lambda_1}} \xi_1 - \sqrt{\frac{\lambda_1}{\lambda_n}} \xi_n \right] = c \rho \left[\sqrt{\frac{\lambda_n}{\lambda_1}} \xi_1 - \sqrt{\frac{\lambda_1}{\lambda_n}} \xi_n \right].
 \end{aligned}$$

Поскольку

$$\|z^0\|_A^2 = (Az^0, z^0) = c^2 \sqrt{\lambda_1 \lambda_n} \left(\sqrt{\frac{\lambda_n}{\lambda_1}} + \sqrt{\frac{\lambda_1}{\lambda_n}} \right) = c^2 (\lambda_1 + \lambda_n),$$

а

$$\begin{aligned}
 Az^1 &= c \rho \left(\sqrt{\lambda_1 \lambda_n} \xi_1 - \sqrt{\lambda_1 \lambda_n} \xi_n \right), \\
 \|z^1\|_A^2 &= c^2 \rho^2 (\lambda_n + \lambda_1) = \rho^2 \|z^0\|_A^2,
 \end{aligned}$$

то

$$\|z^1\|_A = \rho \|z^0\|_A.$$

Делая следующую итерацию, найдем, что

$$z^2 = \rho^2 z^0$$

и т.д. Отсюда вытекает, что

$$\|z^k\|_A = \rho^k \|z^0\|_A,$$

т.е. полученная оценка точная.

Следует, однако, заметить, что такие плохие начальные приближения в реальных задачах практически не встречаются, и итерации, особенно на начальном этапе, сходятся много быстрее. По мере увеличения числа итераций скорость сходимости уменьшается и выходит на ту, которая гарантируется оценкой. Имея хорошее начальное приближение, можно получить приближенное решение с хорошей точностью при существенно меньших трудозатратах.

§ 10

Метод сопряженных градиентов

Построим другой метод минимизации функции

$$J(x) = \frac{1}{2}(Ax, x) - (b, x), \quad (10.1)$$

точка минимума которой совпадает с решением системы

$$Ax = b, \quad A = A^T > 0. \quad (10.2)$$

В методе скорейшего спуска на каждом шаге происходила одномерная минимизация вдоль направления, задаваемого антиградиентом, который совпадает с невязкой $r^k = b - Ax^k$. Рассмотрим теперь последовательную минимизацию $J(x)$ вдоль совокупности направлений $\{p^1, p^2, \dots\}$, которые не обязаны совпадать с направлениями невязок $\{r^0, r^1, \dots\}$.

Пусть направления p^1, p^2, \dots заданы, и (ср. с (9.7))

$$x^{k+1} = x^k + \alpha_{k+1}p^{k+1}, \quad k = 0, 1, \dots \quad (10.3)$$

Поскольку

$$\begin{aligned} J(x^{k+1}) &= J(x^k + \alpha p^{k+1}) = \\ &= \frac{1}{2} (A(x^k + \alpha p^{k+1}), x^k + \alpha p^{k+1}) - (b, x^k + \alpha p^{k+1}) = \\ &= \frac{1}{2} A(x^k, x^k) + \alpha (Ax^k, p^{k+1}) + \frac{\alpha^2}{2} (Ap^{k+1}, p^{k+1}) - (b, x^k) - \alpha (b, p^{k+1}) = \\ &= J(x^k) + \alpha [(Ax^k, p^{k+1}) - (b, p^{k+1})] + \frac{\alpha^2}{2} (Ap^{k+1}, p^{k+1}) = \\ &= J(x^k) - \alpha (r^k, p^{k+1}) + \frac{\alpha^2}{2} (Ap^{k+1}, p^{k+1}), \end{aligned} \quad (10.4)$$

то, дифференцируя это выражение по α и приравнявая производную нулю, находим итерационный параметр

$$\alpha_{k+1} = \frac{(r^k, p^{k+1})}{(p^{k+1}, Ap^{k+1})}. \quad (10.5)$$

Подставляя (10.5) в (10.4), найдем, что

$$\begin{aligned} J(x^{k+1}) &= J(x^k) - \\ &- \frac{(r^k, p^{k+1})}{(p^{k+1}, Ap^{k+1})} (r^k, p^{k+1}) + \frac{1}{2} \frac{(r^k, p^{k+1})^2}{(p^{k+1}, Ap^{k+1})^2} (Ap^{k+1}, p^{k+1}) = \\ &= J(x^k) - \frac{1}{2} \frac{(r^k, p^{k+1})^2}{(p^{k+1}, Ap^{k+1})}, \end{aligned} \quad (10.6)$$

т.е. на $(k+1)$ итерации действительно будет происходить уменьшение функции $J(x)$, если выполнено условие

$$(r^k, p^{k+1}) \neq 0. \quad (10.7)$$

Замечание 10.1. Без ограничения общности можно предполагать, что

$$x^0 = 0. \quad (10.8)$$

Если бы нам было известно хорошее приближение \tilde{x} , то, делая замену $x = \tilde{x} + z$, мы бы нашли, что $Az + A\tilde{x} = b$ и $Az = b - A\tilde{x}$. Тем самым, для z начальным приближением было бы $z^0 = 0$.

Из (10.3) следует, что при начальном приближении (10.8) векторы x^k являются линейными комбинациями векторов p^1, p^2, \dots, p^k , т.е.

$$x^k \in \text{span}\{p^1, p^2, \dots, p^k\}. \quad (10.9)$$

При выборе направлений p^i наша задача состоит в том, чтобы гарантировать сходимость и добиться скорости сходимости бóльшей, чем у метода скорейшего спуска. Представляется, что наилучшим способом выбора p^i был бы такой, при котором x^{k+1} минимизировал бы функцию $J(x)$ не только по направлению p^{k+1} , но и по всему подпространству $\text{span}\{p^1, p^2, \dots, p^{k+1}\} \subset \mathbb{R}^n$, т.е.

$$J(x^{k+1}) = \min_{x \in \text{span}\{p^1, p^2, \dots, p^{k+1}\}} J(x). \quad (10.10)$$

Если бы такой выбор p^i удалось осуществить, то это не только гарантировало бы сходимость, но привело бы к конечности итерационного процесса, ибо при $k+1 = n$ и линейно независимых p^i задача (10.10) представляет собой исходную задачу глобальной минимизации, и, следовательно, $Ax^n = b$.

Попытаемся решить поставленную задачу. Пусть

$$P_k = [p^1 p^2 \dots p^k]$$

есть $(n \times k)$ -матрица, столбцами которой являются искомые направления. Пусть $x = P_k y + \alpha p^{k+1} \in \text{im } P_{k+1}$,¹ $y \in \mathbb{R}^k$, $\alpha \in \mathbb{R}$. Тогда (см. (10.4))

$$J(x) = J(P_k y) + \alpha (AP_k y, p^{k+1}) + \frac{\alpha^2}{2} (Ap^{k+1}, p^{k+1}) - \alpha (b, p^{k+1}). \quad (10.11)$$

¹Напомним, что множество всех векторов x , представимых в виде $x = By$, называется образом матрицы B и обозначается $\text{im } B$.

Если бы в (10.11) отсутствовал "перекрестный" член

$$\alpha(P_k y, Ap^{k+1}),$$

то задача минимизации $J(x)$ на $\text{span}\{p^1, p^2, \dots, p^{k+1}\} = \text{im } P_{k+1}$, т.е. задача (10.10), распалась бы на минимизацию по $\text{im } P_k$, где решение x^k предполагается известным, и простую минимизацию для определения скалярной величины α .

В самом деле, пусть выполнены условия

$$(p^i, Ap^j) = 0, \quad i \neq j. \quad (10.12)$$

(Векторы, удовлетворяющие условию (10.12), называются A -сопряженными или A -ортогональными.) Определим вектор $x^k \in \text{im } P_k$ и $\alpha_{k+1} \in \mathbb{R}$ следующим образом

$$J(x^k) = \min_y J(P_k y), \quad \alpha_{k+1} = \frac{(b, p^{k+1})}{(p^{k+1}, Ap^{k+1})}. \quad (10.13)$$

Тогда (см. (10.11))

$$\min_{y, \alpha} J(P_k y + \alpha p^{k+1}) = \min_y J(P_k y) + \min_{\alpha} \left\{ \frac{\alpha^2}{2} (Ap^{k+1}, p^{k+1}) - \alpha (b, p^{k+1}) \right\}$$

находится при $P_k y = x^k$ и $\alpha = \alpha_{k+1}$ из (10.13). Покажем, что на самом деле α_{k+1} из (10.13) совпадает с (10.5). В силу (10.9) и (10.12)

$$(Ap^{k+1}, x^k) = 0$$

и, следовательно,

$$(p^{k+1}, b) = (p^{k+1}, b - Ax^k + Ax^k) = (p^{k+1}, r^k), \quad (10.14)$$

что вместе с (10.13) приводит к (10.5).

Итак, для реализации задуманного метода нужно последовательно находить A -сопряженные векторы p^1, p^2, \dots, p^{k+1} , для которых выполнено условие (10.7), и проводить вычисления по формуле (10.3) с параметром α_{k+1} из (10.5).

Обратимся к наиболее целесообразному выбору векторов p^{k+1} . При выборе p^{k+1} наша цель состоит в быстрой минимизации функции $J(x)$, и в силу (10.6) мы должны максимизировать

$$\frac{(r^k, p^{k+1})^2}{(p^{k+1}, Ap^{k+1})}.$$

Замечание 10.2. Эта величина не зависит от длины вектора p^{k+1} , а зависит только от его направления. Поэтому при отыскании p^{k+1} достаточно ограничиться нахождением его направления.

Поскольку p^{k+1} должен еще удовлетворять условиям A -сопряженности (10.12), т.е. быть ортогональным к $\{Ap^1, Ap^2, \dots, Ap^k\}$, то искомым вектор

$$p^{k+1} \in (\text{span} \{Ap^1, Ap^2, \dots, Ap^k\})^\perp = (\text{im } AP_k)^\perp.$$

Пусть $r^k = r_{\parallel}^k + r_{\perp}^k$, где $r_{\parallel}^k \in \text{im } (AP_k)$, а $r_{\perp}^k \in (\text{im } (AP_k))^\perp$. Тогда

$$(r^k, p^{k+1}) = (r_{\parallel}^k + r_{\perp}^k, p^{k+1}) = (r_{\perp}^k, p^{k+1}) = \|r_{\perp}^k\| \|p^{k+1}\| \cos(r_{\perp}^k, p^{k+1})$$

и искомым максимум будет достигаться при $|\cos(r_{\perp}^k, p^{k+1})| = 1$, т.е., например, при

$$p^{k+1} = r_{\perp}^k \in (\text{im } (AP_k))^\perp \quad (10.15)$$

— ортогональной проекции r^k на $(\text{im } (AP_k))^\perp$. Отметим, что отсюда следует соотношение

$$p^1 = r^0. \quad (10.16)$$

Построение процесса минимизации $J(x)$ в первом приближении будет закончено, если принять во внимание, что имеет место

Теорема 10.1. *Два последовательных направления спуска в методе сопряженных градиентов связаны соотношением*

$$p^{k+1} = r^k + \beta_{k+1} p^k. \quad (10.17)$$

Доказательство этой теоремы мы отложим на потом, а сейчас заметим, что поскольку векторы p^k и p^{k+1} должны быть A -сопряжены, то для параметра β_{k+1} из (10.17) имеет место представление

$$\beta_{k+1} = -\frac{(r^k, Ap^k)}{(p^k, Ap^k)}. \quad (10.18)$$

Итак, метод сопряженных градиентов состоит в вычислениях по следующим формулам

$$\begin{aligned} r^k &= b - Ax^k, & k &= 0, 1, \dots, \\ p^{k+1} &= r^k + \beta_{k+1} p^k, & k &= 1, 2, \dots, & p^1 &= r^0, \\ x^{k+1} &= x^k + \alpha_{k+1} p^{k+1}, & k &= 0, 1, \dots, & x^0 &= 0, \\ \alpha_{k+1} &= (r^k, p^{k+1}) / (p^{k+1}, Ap^{k+1}), & k &= 0, 1, \dots, \\ \beta_{k+1} &= -(Ap^k, r^k) / (Ap^k, p^k), & k &= 1, 2, \dots, \end{aligned} \quad (10.19)$$

Дадим оценку скорости сходимости метода сопряженных градиентов. Имеет место

Теорема 10.2. *Метод сопряженных градиентов (10.19) сходится не хуже, чем чебышевский итерационный метод, т.е.*

$$\|x^k - x\|_A \leq 2 \left[(1 - \sqrt{\lambda_1/\lambda_n}) / (1 + \sqrt{\lambda_1/\lambda_n}) \right]^k \|x\|_A,$$

где λ_1 и λ_n — минимальное и максимальное собственные значения матрицы A .

Доказательство. Покажем сначала, что минимизация $J(x^k)$ ведет к минимизации $\|x^k - x\|_A$. В самом деле, пусть

$$z^k = x^k - x.$$

Тогда, подставляя $x^k = x + z^k$ в $J(x^k)$ и принимая во внимание (10.4) с заменой x^k на x , $\alpha_{k+1}p^{k+1}$ на z^k , а x^{k+1} на x^k , будем иметь

$$J(x^k) = \frac{1}{2}\|z^k\|_A^2 + J(x). \quad (10.20)$$

Установим теперь связь между z^k на последовательных итерациях. Из третьего соотношения (10.19) находим, что

$$p^{k+1} = (x^{k+1} - x^k)/\alpha_{k+1}.$$

Подставим это представление p^{k+1} во второе соотношение (10.19)

$$\frac{x^{k+1} - x^k}{\alpha_{k+1}} - \beta_{k+1} \frac{x^k - x^{k-1}}{\alpha_k} = b - Ax^k.$$

Отсюда находим, что

$$\frac{z^{k+1} - z^k}{\alpha_{k+1}} - \beta_{k+1} \frac{z^k - z^{k-1}}{\alpha_k} + Az^k = 0.$$

Далее,

$$z^1 = z^0 + \alpha_1 p^1 = z^0 + \alpha_1 r^0 = z^0 + \alpha_1 b = z^0 + \alpha_1 Ax = z^0 - \alpha_1 Az^0.$$

Тем самым,

$$z^k = p_k(A)z^0, \quad p_k(0) = 1$$

и

$$\|z^k\|_A = \|p_k(A)z^0\|_A.$$

Но по построению x^k и с учетом (10.20)

$$\|z^k\|_A = \min_{q_k} \|q_k(A)z^0\|_A, \quad q(0) = 1$$

и, следовательно,

$$\begin{aligned} \|z^k\|_A &\leq \|q_k(A)z^0\|_A = \|q_k(A)A^{1/2}z^0\|_2 \leq \|q_k(A)\| \|z^0\|_A \leq \max_{\lambda_1 \leq t \leq \lambda_n} |q_k(t)| \|z^0\|_A = \\ &= \max_{y \in [-1,1]} \left| q_k \left(\frac{\lambda_n + \lambda_1}{2} - \frac{\lambda_n - \lambda_1}{2} y \right) \right| \|z^0\|_A = \max_{y \in [-1,1]} |\hat{q}_k(y)| \|z^0\|_A. \end{aligned}$$

Если положить $\hat{Q}_k(y) = q_k T_k(y)$ (см. лекцию 8), то

$$\|z^k\|_A \leq q_k \|z^0\|_A = \frac{2\rho_1^k}{1 + 2\rho_1^{2k}} \|z^0\|_A \leq 2 \left(\frac{1 - \sqrt{\lambda_1/\lambda_n}}{1 + \sqrt{\lambda_1/\lambda_n}} \right)^k \|z^0\|_A.$$

Теорема доказана.

Чтобы описание метода сопряженных градиентов (10.19) было корректным, нужно доказать теорему 10.1 (о связи между двумя последовательными направлениями спуска). Для этого нам понадобится ряд вспомогательных утверждений.

Лемма 10.1. Пусть p^1, p^2, \dots, p^k суть ненулевые A -сопряженные векторы. Тогда, либо существует A -сопряженный к ним вектор p^{k+1} , удовлетворяющий условию (10.7), либо $r^k = 0$.

Доказательство. Пусть не существует такого A -сопряженного с p^1, p^2, \dots, p^k вектора p^{k+1} , для которого выполнено условие (10.7), т.е. для любого вектора

$$p \perp \{Ap^1, Ap^2, \dots, Ap^k\} = \text{im } AP_k \quad (10.21)$$

имеет место равенство

$$(r^k, p) = 0. \quad (10.22)$$

Напомним, что в силу (10.14) для любого вектора p , удовлетворяющего (10.21),

$$(p, r^k) = (p, b)$$

и, следовательно, (см. (10.22)),

$$(p, Ax) = 0.$$

Тем самым, решение $x \in \text{im } P_k$. Но x^k минимизирует $J(x)$ на $\text{im } P_k$ и, следовательно, $x^k = x$, т.е. $r^k = 0$. Отсюда же вытекает, что, если $r^k \neq 0$, то существует p^{k+1} из (10.21), для которого выполнено условие (10.7). Лемма доказана.

Замечание 10.3. Лемма 10.1 утверждает, что либо мы на k -ой итерации закончили вычисления, получив точное решение задачи (10.2), либо имеем возможность вычисления продолжить.

Теорема 10.3. После k итераций метода сопряженных градиентов при каждом j от 1 до k

$$\begin{aligned} \text{span } \{p^1, p^2, \dots, p^j\} &= \text{span } \{r^0, r^1, \dots, r^{j-1}\} = \\ &= \text{span } \{b, Ab, \dots, A^{j-1}b\} =: \mathcal{K}_j(A, b). \end{aligned} \quad (10.23)$$

Доказательство проведем методом полной математической индукции. При $j = 1$ соотношения (10.23) имеют место, ибо в силу выбора (10.8) начального приближения $r^0 = b$, а из (10.16) $p^1 = r^0$. Предположим, что (10.23) справедливы при некотором j , удовлетворяющем неравенству $1 \leq j < k$. Докажем их справедливость при $j + 1$.

В качестве первого шага покажем, что

$$\text{span } \{p^1, p^2, \dots, p^{j+1}\} \subset \text{span } \{r^0, r^1, \dots, r^j\}. \quad (10.24)$$

В силу (10.15)

$$p^{j+1} = r_{\perp}^j = r^j - r_{\parallel}^j, \quad r_{\parallel}^j \in \text{im } [AP_j]$$

и, следовательно,

$$p^{j+1} = r^j - AP_j y_j, \quad y_j \in \mathbb{R}^j. \quad (10.25)$$

Из (10.3) вытекает, что

$$Ax^j = Ax^{j-1} + \alpha_j Ap^j.$$

Вычитая из обеих частей этого равенства по b , будем иметь

$$r^j = r^{j-1} - \alpha_j Ap^j \quad (10.26)$$

и, следовательно,

$$Ap^j = -(r^j - r^{j-1})/\alpha_j.$$

Подставляя это представление в (10.25), находим, что

$$p^{j+1} = r^j + \left[\frac{r^1 - r^0}{\alpha_1} \frac{r^2 - r^1}{\alpha_2} \dots \frac{r^j - r^{j-1}}{\alpha_j} \right] y_j.$$

Отсюда и из предположения индукции следует включение (10.24).

Теперь установим включение

$$\text{span} \{r^0, r^1, \dots, r^j\} \subset \text{span} \{b, Ab, \dots, A^j b\}. \quad (10.27)$$

По предположению индукции векторы $p^j \in \text{span} \{b, Ab, \dots, A^{j-1} b\}$. Поэтому

$$Ap^j \in \text{span} \{b, Ab, \dots, A^j b\}.$$

Если теперь принять во внимание включение

$r^{j-1} \in \text{span} \{b, Ab, \dots, A^{j-1} b\}$, то из (10.26) найдем, что

$$r^j \in \text{span} \{b, Ab, \dots, A^j b\}.$$

Вновь принимая во внимание предположение индукции, будем иметь желаемое включение (10.27).

Итак, вместо равенства (10.23) мы пока имеем только включения (10.24), (10.27) пространств, размерность каждого из которых не превышает $j+1$. Поскольку векторы p^1, p^2, \dots, p^{j+1} ненулевые и A -сопряженные, то

$$\dim \text{span} \{p^1, p^2, \dots, p^{j+1}\} = j + 1.$$

Отсюда и из включений (10.24), (10.27) следует искомое равенство (10.23).

Лемма 10.2. После k итераций по методу сопряженных градиентов невязка r^j ортогональна всем векторам спуска p^1, \dots, p^j , т.е.

$$P_j^T r^j = 0, \quad j = \overline{1, k}. \quad (10.28)$$

Доказательство. В силу (10.9) существует вектор $y_j \in \mathbb{R}^j$ такой, что $x^j = P_j y_j$. Поэтому

$$\begin{aligned} J(x^j) &= J(P_j y_j) = \frac{1}{2}(AP_j y_j, P_j y_j)_n - (b, P_j y_j)_n = \\ &= \frac{1}{2}[P_j y_j]^T [AP_j y_j] - [P_j y_j]^T b = \\ &= \frac{1}{2}y_j^T P_j^T AP_j y_j - y_j^T P_j^T b = \\ &= \frac{1}{2}([P_j^T AP_j]y_j, y_j)_j - (P_j^T b, y_j)_j, \end{aligned}$$

т.е. y_j есть решение задачи минимизации с матрицей $P_j^T AP_j$ и вектором $P_j^T b$, и, следовательно, вектор y_j является решением следующей системы

$$[P_j^T AP_j]y_j = P_j^T b.$$

Отсюда

$$P_j^T r^j = P_j^T (b - Ax^j) = P_j^T (b - AP_j y_j) = 0.$$

Лемма доказана.

Теорема 10.4. После k шагов метода сопряженных градиентов невязки r^0, r^1, \dots, r^k взаимно ортогональны.

Доказательство. В силу теоремы 10.3

$$p^j \in \text{span} \{r^0, r^1, \dots, r^{j-1}\}.$$

Это означает, что p^1 выражается только через r^0 , p^2 — только через r^0 и r^1 и т.д., т.е.

$$P_j = [p^1 p^2 \dots p^j] = [r^0 r^1 \dots r^{j-1}] U_j =: R_j U_j, \quad (10.29)$$

где U_j — верхняя треугольная $(j \times j)$ -матрица.

Поскольку векторы p^1, p^2, \dots, p^j , а в силу теоремы 10.3 и векторы r^0, r^1, \dots, r^{j-1} линейно независимы, то U_j — невырожденная матрица. Подставляя представление (10.29) матрицы P_j в (10.28), будем иметь

$$0 = U_j^T R_j^T r^j = U_j^T [(r^0, r^j) (r^1, r^j) \dots (r^{j-1}, r^j)]^T.$$

Рассматривая это соотношение как систему линейных однородных уравнений относительно (r^i, r^j) с невырожденной матрицей, приходим к заключению, что $(r^i, r^j) = 0$ при $i \neq j$. Теорема доказана.

Мы теперь имеем все необходимое для того, чтобы доказать теорему 10.1, т.е.

$$p^{k+1} = r^k + \beta_{k+1} p^k. \quad (10.30)$$

Доказательство теоремы 10.1. В силу (10.15)

$$p^{k+1} = r_{\perp}^k = r^k - r_{\parallel}^k, \quad r_{\parallel}^k \in \text{im } AP^k,$$

т.е.

$$p^{k+1} = r^k - \sum_{j=1}^k c_{k+1,j} Ap^j.$$

Поскольку в силу теоремы 10.3 $p^j \in K_j(A, b)$, то

$$Ap^j \in K_{j+1}(A, b), \quad (10.31)$$

а, снова используя теорему 10.3 находим, что $Ap^j \in \text{span} \{p^1, p^2, \dots, p^{j+1}\}$ и, следовательно,

$$p^{k+1} = r^k - \sum_{j=1}^{k+1} d_{k+1,j} p^j$$

или

$$(1 + d_{k+1,k+1})p^{k+1} = r^k - \sum_{j=1}^k d_{k+1,j} p^j. \quad (10.32)$$

Коэффициент $(1 + d_{k+1,k+1})$ при p^{k+1} нулю не равен, ибо в противном случае

$$r^k = \sum_{j=1}^k d_{k+1,j} p^j = P_k [d_{k+1,1} d_{k+1,2} \dots d_{k+1,k}]^T = P_k d^k$$

и с учетом леммы 10.2

$$\|r^k\|^2 = r^{kT} r^k = r^{kT} P_k d^k = (r^{kT} P_k d^k)^T = d^{kT} (P_k^T r^k) = 0,$$

т.е. $r^k = 0$, что в силу леммы 10.1 возможно лишь по завершении итераций.

Так как вектор p^{k+1} из (10.32) должен быть A -ортогонален векторам p^j , $j = 1, 2, \dots, k$, то

$$d_{k+1,j} = \frac{(r^k, Ap^j)}{(Ap^j, p^j)}, \quad j = 1, 2, \dots, k. \quad (10.33)$$

Снова обращаясь к (10.31) и теореме 10.3, находим, что

$$Ap^j \in \text{span} \{r^0, r^1, \dots, r^j\},$$

а по теореме 10.4

$$(r^k, Ap^j) = \sum_{i=0}^j l_{j,i} (r^k, r^i) = 0, \quad j = 1, 2, \dots, k-1.$$

Следовательно,

$$d_{k+1,j} = 0, \quad j = 1, 2, \dots, k-1,$$

а (10.32) принимает вид

$$(1 + d_{k+1,k+1})p^{k+1} = r^k - d_{k+1,k}p^k,$$

где $d_{k+1,k}$ определяется соотношением (10.33) (ср. с (10.18)), что с точностью до длины вектора p^{k+1} (см. Замечание 10.2) совпадает с (10.17). Теорема доказана.

Преобразуем соотношения (10.19) метода сопряженных градиентов. В этих соотношениях наиболее трудоемкими являются две операции: вычисление векторов Ax^k и Ap^k . Однако операцию вычисления вектора Ax^k можно исключить. Поскольку этот вектор используется только при вычислении невязки r^k , то можно заменить первую из формул (10.19) на (10.26)

$$r^k = r^{k-1} - \alpha_k Ap^k, \quad k = 1, 2, \dots, \quad r^0 = b. \quad (10.34)$$

Преобразуем еще формулы для вычисления параметров α_{k+1} и β_{k+1} . Подставляя второе из соотношений (10.19) в четвертое и принимая во внимание лемму 10.2, найдем, что

$$\alpha_{k+1} = (r^k, r^k) / (p^{k+1}, Ap^{k+1}), \quad k = 0, 1, \dots \quad (10.35)$$

Далее, заменяя здесь $k+1$ на k и подставляя полученное выражение для (p^k, Ap^k) в последнее из соотношений (10.19), будем иметь

$$\beta_{k+1} = -\alpha_k \frac{(Ap^k, r^k)}{(r^{k-1}, r^{k-1})}.$$

Теперь подставим сюда вместо Ap^k его выражение из (10.34). Принимая во внимание теорему 10.4, найдем, что

$$\beta_{k+1} = \frac{(r^k, r^k)}{(r^{k-1}, r^{k-1})}, \quad k = 1, 2, \dots \quad (10.36)$$

С учетом (10.34)-(10.36) формулы метода сопряженных градиентов (10.19) преобразуются к виду

$$\begin{aligned} r^k &= r^{k-1} - \alpha_k Ap^k, \quad k = 1, 2, \dots, \quad r^0 = b, \\ p^{k+1} &= r^k + \beta_{k+1} p^k, \quad k = 1, 2, \dots, \quad p^1 = r^0, \\ x^{k+1} &= x^k + \alpha_{k+1} p^{k+1}, \quad k = 0, 1, \dots, \quad x^0 = 0, \\ \alpha_{k+1} &= \|r^k\|^2 / (p^{k+1}, Ap^{k+1}), \quad k = 0, 1, \dots, \\ \beta_{k+1} &= \|r^k\|^2 / \|r^{k-1}\|^2, \quad k = 1, 2, \dots \end{aligned} \quad (10.37)$$

Легко проверить, что вычисления можно проводить в следующем порядке

$$\begin{aligned} r^0 &= b, \quad p^1 = r^0, \quad Ap^1, \quad \alpha_1, \quad x^1, \\ r^1, \quad \beta_2, \quad p^2, \quad Ap^2, \quad \alpha_2, \quad x^2, \quad \dots \end{aligned}$$

§ 11

Проблема собственных значений

11.1 Постановка задачи

Пусть A - квадратная матрица с действительными коэффициентами, и требуется найти собственные векторы и собственные значения этой матрицы. Напомним

Определение 11.1. Число λ называется собственным значением матрицы A , если однородная система

$$A\xi = \lambda\xi \quad (11.1)$$

имеет нетривиальное решение $\|\xi\| \neq 0$. Это нетривиальное решение называется собственным вектором матрицы A , отвечающим собственному значению λ .

Собственные значения являются нулями характеристического многочлена

$$\det [A - \lambda I] = 0,$$

степень которого совпадает с порядком матрицы и есть n . Тем самым, у каждой квадратной матрицы существует n собственных значений, действительных или комплексных, простых или кратных. С собственными векторами ситуация сложнее: их число может быть от 1 до n .

Определение 11.2. Матрицы A и B называются подобными, если существует невырожденная матрица S (матрица подобия) такая, что $B = S^{-1}AS$.

Подобные матрицы имеют одинаковый набор собственных значений.

Любая матрица A преобразованием подобия $S^{-1}AS$ с подходящей матрицей подобия S может быть приведена к нормальной (Жордановой) форме. (На главной диагонали — собственные значения, а на наддиагонали — нули и (или) единицы)

$$\begin{bmatrix} \lambda_1 & \sigma_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & \lambda_2 & \sigma_2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \lambda_3 & \sigma_3 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & \lambda_{n-1} & \sigma_{n-1} \\ 0 & 0 & 0 & 0 & \dots & 0 & \lambda_n \end{bmatrix}$$

Определение 11.3. Матрица A называется матрицей простой структуры (или диагонализуемой), если ее жордановой формой является диагональная матрица.

- Матрица простой структуры имеет ровно n линейно независимых собственных векторов. Про такую матрицу еще говорят, что она имеет полный набор собственных векторов.
- Если все собственные значения матрицы A различны, то она заведомо имеет простую структуру.
- Симметричная матрица имеет простую структуру, и поэтому у нее имеется n линейно-независимых собственных векторов. Ее собственные векторы, отвечающие различным собственным значениям, ортогональны в смысле обычного скалярного произведения

$$y^T x = (x, y) = \sum_{i=1}^n x_i y_i,$$

а собственные векторы, отвечающие кратному собственному значению (собственному значению кратности m отвечает m линейно-независимых собственных векторов), могут быть ортогонализированы.

- Матрица A^T имеет те же собственные значения, что и матрица A , а собственные векторы ξ_i и η_j матриц A и A^T , соответственно, отвечающие различным собственным значениям, ортогональны (образуют ортонормированную систему).
- Собственные векторы матриц A и A^{-1} совпадают, а собственные значения связаны соотношением $\lambda_i(A^{-1}) = \lambda_i^{-1}(A)$.
- Собственные векторы матриц A и $B = A + \alpha I$ совпадают, а собственные значения связаны соотношениями $\lambda_i(B) = \lambda_i(A) + \alpha$.

Задача нахождения всех собственных значений и собственных векторов называется полной проблемой собственных значений. Эта проблема в общем случае довольно сложна.

Наряду с полной проблемой собственных значений существуют частичные проблемы собственных значений, отыскание решений которых много проще.

К последним относятся:

- 1) Задача отыскания максимального или минимального по модулю собственного значения и, быть может, отвечающего ему собственного вектора.
- 2) Задача отыскания двух наибольших по модулю собственных значений и соответствующих собственных векторов.
- 3) Задача отыскания собственного значения, наиболее близкого к заданному числу.

Этими задачами мы и займемся.

11.2 Степенной метод решения частных проблем

Изложим метод, позволяющий решить некоторые из частных проблем собственных значений при помощи вычислений последовательных итераций произвольного вектора. Излагаемый метод называется степенным и является простейшим итерационным методом.

11.2.1 Нахождение максимального по модулю собственного значения

Будем предполагать, что матрица A имеет простую структуру, а ее собственные значения действительны, т.е.

$$A\xi_i = \lambda_i\xi_i, \quad \text{Im } \lambda_i = 0, \quad i = \overline{1, n}, \quad (11.2)$$

$$\|\xi_i\| = \|\xi_i\|_2 = (\xi_i, \xi_i)^{1/2} = 1, \quad i = \overline{1, n}. \quad (11.3)$$

Допустим, что

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|. \quad (11.4)$$

Зададим произвольный вектор x^0 . Его разложение по собственным векторам ξ_i матрицы A имеет вид

$$x^0 = c_1\xi_1 + c_2\xi_2 + \dots + c_n\xi_n. \quad (11.5)$$

Здесь c_1, c_2, \dots, c_n — координаты вектора x^0 в базисе $\xi_1, \xi_2, \dots, \xi_n$. Предположим, что

$$c_1 \neq 0 \quad (11.6)$$

и вычислим последовательно векторы

$$x^k = Ax^{k-1}, \quad k = 1, 2, \dots \quad (11.7)$$

Тогда согласно (11.5), (11.2)

$$\begin{aligned} x^1 &= Ax^0 = A(c_1\xi_1 + c_2\xi_2 + \dots + c_n\xi_n) = \\ &= c_1\lambda_1\xi_1 + c_2\lambda_2\xi_2 + \dots + c_n\lambda_n\xi_n \end{aligned}$$

и вообще

$$x^k = c_1 \lambda_1^k \xi_1 + c_2 \lambda_2^k \xi_2 + \cdots + c_n \lambda_n^k \xi_n = \lambda_1^k (c_1 \xi_1 + \eta^k), \quad (11.8)$$

где

$$\eta^k = c_2 (\lambda_2/\lambda_1)^k \xi_2 + c_3 (\lambda_3/\lambda_1)^k \xi_3 + \cdots + c_n (\lambda_n/\lambda_1)^k \xi_n.$$

Вычисляя норму η^k , с учетом (11.3), (11.4) находим, что

$$\begin{aligned} \|\eta^k\|_2 &\leq \sum_{j=2}^n |c_j| |\lambda_j/\lambda_1|^k \|\xi_j\|_2 \leq |\lambda_2/\lambda_1|^k \sum_{j=2}^n |c_j| = \\ &= O(|\lambda_2/\lambda_1|^k) \rightarrow 0 \quad \text{при } k \rightarrow \infty. \end{aligned} \quad (11.9)$$

С учетом (11.8) вычислим скалярные произведения

$$\begin{aligned} (x^k, x^{k+1}) &= \lambda_1^{2k+1} (c_1 \xi_1 + \eta^k, c_1 \xi_1 + \eta^{k+1}) = \\ &= \lambda_1^{2k+1} [c_1^2 (\xi_1, \xi_1) + c_1 (\xi_1, \eta^{k+1}) + c_1 (\eta^k, \xi_1) + (\eta^k, \eta^{k+1})]. \end{aligned}$$

Оценивая, находим, что

$$\begin{aligned} |(\xi_1, \eta^{k+1})| &\leq \|\xi_1\| \|\eta^{k+1}\| = \|\eta^{k+1}\| \\ |(\eta^k, \xi_1)| &\leq \|\eta^k\|, \quad |(\eta^k, \eta^{k+1})| \leq \|\eta^k\| \|\eta^{k+1}\| \end{aligned}$$

и с учетом (11.9)

$$(x^k, x^{k+1}) = \lambda_1^{2k+1} (c_1^2 + O(|\lambda_2/\lambda_1|^k)).$$

Аналогично

$$(x^k, x^k) = \lambda_1^{2k} (c_1^2 + O(|\lambda_2/\lambda_1|^k)) \quad (11.10)$$

и, следовательно,

$$\lambda_1^{(k)} := \frac{(x^{k+1}, x^k)}{(x^k, x^k)} = \lambda_1 + O(|\lambda_2/\lambda_1|^{k-1}). \quad (11.11)$$

Из (11.10)

$$\|x^k\| = |\lambda_1|^k (|c_1| + O(|\lambda_2/\lambda_1|^k)),$$

а с учетом (11.8)

$$\xi^k := \frac{x^k}{\|x^k\|} = \pm \xi_1 + r^k, \quad (11.12)$$

где

$$\|r^k\| = O(|\lambda_2/\lambda_1|^k).$$

Таким образом, итерационный процесс (11.7), (11.11), (11.12) позволяет найти с любой точностью однократное максимальное по модулю собственное значение (11.11) и отвечающий ему собственный вектор (11.12), если выполнено условие (11.6).

Замечание 11.1. Если $|\lambda_1| > 1$, то $\|x^k\| \rightarrow \infty$ при $k \rightarrow \infty$, а если $|\lambda_1| < 1$, то $\|x^k\| \rightarrow 0$. И то, и другое явление нежелательны при вычислениях на компьютере. В первом случае может произойти переполнение, а во втором случае x^k может стать машинным нулем. Поэтому вместо (11.7) итерации нужно вести по формулам

$$\begin{aligned} \xi_1^0 &= x^0 / \|x^0\|, & x^{k+1} &= A\xi_1^k, \\ \lambda_1^{(k)} &= (x^{k+1}, \xi_1^k) = (A\xi_1^k, \xi_1^k), & \xi_1^{k+1} &= x^{k+1} / \|x^{k+1}\|. \end{aligned} \quad (11.13)$$

Замечание 11.2. Если условие (11.6) не выполнено (априори проверить это условие нельзя), то это еще не значит, что итерационный процесс (11.7) (или (11.13)) с начальным приближением (11.5) не приведет к результату. При достаточно большом числе итераций за счет ошибок округления может появиться ненулевая компонента c_1 , и итерационный процесс выйдет в конце концов на нужное решение. Но при этом нужно иметь в виду, что если $|\lambda_3| \ll |\lambda_2|$, то итерации очень быстро выйдут на второе собственное значение и второй собственный вектор, и можно обмануться, приняв их за искомые величины. Это не так вероятно, если $|\lambda_2|$ и $|\lambda_3|$ не слишком сильно различаются, а требуемая точность достаточно велика. Итерации в этом случае будут сходиться достаточно медленно, и их потребуется много для получения требуемой точности. За это время погрешности округления накопятся, и может сформироваться новая точка притяжения итерационного процесса — (λ_1, ξ_1) . Если нет уверенности в правильности найденного собственного значения, следует провести еще один или несколько расчетов с другими начальными приближениями.

Замечание 11.3. 1) Подтверждением того, что λ_1 не является кратным собственным значением, и что нет собственного значения $(-\lambda_1)$, служит сходимость итерационного процесса к одному и тому же собственному вектору (с точностью до знака) при различных начальных приближениях.

2) Если при различных начальных векторах x^0 значения $\lambda_1^{(k)}$ сходятся к одному и тому же числу, а последовательности векторов $\xi_1^{(k)}$ приводят к неколлинеарным векторам, то это обстоятельство служит подтверждением того, что максимальное по модулю собственное значение является кратным. Если требуется найти собственное подпространство, или нужно определить кратность найденного собственного значения, нужно проводить вычисления с различными начальными приближениями до тех пор, пока перестанут получаться векторы, линейно-независимые с уже найденными.

3) Если значения $\lambda_1^{(k)}$ не сходятся при $k \rightarrow \infty$, однако $\lambda_1^{(2k+1)}$ и $\lambda_1^{(2k)}$ сходятся, но к разным числам, то это свидетельствует о наличии двух максимальных по модулю собственных значений, знаки которых различны. В этом случае целесообразно произвести "сдвиг спектра" путем проведения итераций (11.7) с матрицей $A' = A + cI$, где c — заданное число.

11.2.2 Нахождение второго по величине модуля собственного значения

Пусть

$$|\lambda_1| > |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|.$$

Будем считать, что λ_1 , ξ_1 и η_1 ($A^T \eta_1 = \lambda_1 \eta_1$) известны, причем $\|\xi_1\| = 1$, $(\eta_1, \xi_1) = 1$. Найти λ_1 , ξ_1 и η_1 можно описанным выше способом. Пусть x^0 — произвольный вектор, такой, что $(x^0, \eta_2) \neq 0$. Тогда

$$x^0 = c_1 \xi_1 + c_2 \xi_2 + \dots + c_n \xi_n, \quad c_1 = (x^0, \eta_1), \quad c_2 \neq 0.$$

Построим вектор

$$y^0 = x^0 - (x^0, \eta_1) \xi_1 = c_2 \xi_2 + c_3 \xi_3 + \dots + c_n \xi_n$$

и вектор

$$\xi_2^0 = y^0 / \|y^0\|.$$

Итерационный процесс будем осуществлять по формулам

$$\begin{aligned} x^{k+1} &= A \xi_2^k, & \lambda_2^{(k)} &= (x^{k+1}, \xi_2^k), \\ y^{k+1} &= x^{k+1} - (x^{k+1}, \eta_1) \xi_1, & \xi_2^{k+1} &= y^{k+1} / \|y^{k+1}\|. \end{aligned}$$

Тогда

$$\begin{aligned} \lambda_2^{(k)} &= \lambda_2 + O(|\lambda_3/\lambda_2|^k), \\ \xi_2^k &= \pm \xi_2 + O(|\lambda_3/\lambda_2|^k). \end{aligned}$$

11.2.3 Нахождение $\max_{1 \leq i \leq n} \lambda_i(A)$ и $\min_{1 \leq i \leq n} \lambda_i(A)$

а) Найдем максимальное по модулю собственное значение по описанной выше методике. Пусть это $\bar{\lambda}(A)$

$$|\bar{\lambda}(A)| = \max_i |\lambda_i(A)|.$$

Если $\bar{\lambda}(A) > 0$, то найденное максимальное по модулю собственное значение $\bar{\lambda}(A)$ будет искомым максимальным значением

$$\max_i \lambda_i(A) = \bar{\lambda}(A) \tag{11.14}$$

б) Найдем $\bar{\lambda}_i(B)$: $|\bar{\lambda}_i(B)| = \max_i |\lambda_i(B)|$, где

$$B = A - \bar{\lambda}(A)I.$$

При этом

$$\lambda_i(B) = \lambda_i(A) - \bar{\lambda}(A) \leq 0.$$

Поэтому максимальное по модулю собственное значение матрицы B есть минимальное собственное значение этой матрицы

$$\bar{\lambda}(B) = \min_i [\lambda_i(A) - \bar{\lambda}(A)] = \min_i \lambda_i(A) - \bar{\lambda}(A),$$

т.е.

$$\min_i \lambda_i(A) = \bar{\lambda}(A) + \bar{\lambda}(B). \quad (11.15)$$

Если же $\bar{\lambda}(A) < 0$, то все наоборот.

11.3 Метод обратных итераций

Если матрица A невырождена, то уравнение (11.1) можно переписать в виде

$$A^{-1}\xi = \frac{1}{\lambda}\xi \quad (11.16)$$

и для отыскания наименьшего по модулю собственного значения матрицы A использовать приближение к наибольшему по модулю значению матрицы A^{-1} . Именно, пусть

$$x^0 = c_1\xi_1 + c_2\xi_2 + \dots + c_n\xi_n, \quad c_n \neq 0, \\ |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_{n-1}| > |\lambda_n|.$$

Тогда

$$\frac{1}{|\lambda_n|} > \frac{1}{|\lambda_{n-1}|} \geq \dots \geq \frac{1}{|\lambda_1|}, \\ \xi_n^0 = \frac{x^0}{\|x^0\|}, \quad x^{k+1} = A^{-1}\xi_n^k, \\ \frac{1}{\lambda^{(k+1)}} = (x^{k+1}, \xi_n^k), \quad \xi_n^{k+1} = \frac{x^{k+1}}{\|x^{k+1}\|}, \quad (11.17)$$

Разумеется, вычислять A^{-1} нет необходимости — достаточно решать системы

$$Ax^{k+1} = \xi_n^k \quad (11.18)$$

с одной и той же матрицей и различными правыми частями.

Метод обратных итераций может быть использован и в том случае, когда уже известно с некоторой точностью какое-либо собственное значение, и нужно его уточнить, а также найти отвечающий ему собственный вектор. Для этого нужно вместо матрицы A использовать матрицу

$$A - \tilde{\lambda}I, \quad (11.19)$$

где $\tilde{\lambda}$ — известное приближение к искомому собственному значению. Если приближение достаточно хорошее, то у матрицы $A - \tilde{\lambda}I$ есть собственное значение, по модулю значительно меньшее остальных. В этом случае обратные итерации быстро сходятся к этому малому собственному значению, т.е. к поправке для $\tilde{\lambda}$, а заодно находится (или уточняется) приближенный собственный вектор.

11.4 Пример

Применим степенной метод для отыскания максимального по модулю собственного значения матрицы

$$A = \begin{bmatrix} 2 & 1 \\ -1 & 0 \end{bmatrix}$$

В качестве начального приближения возьмем вектор $x^0 = [0, 1]^T$. Тогда

$$\begin{aligned} x^1 &= \begin{bmatrix} 2 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, & x^2 &= \begin{bmatrix} 2 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \\ x^3 &= \begin{bmatrix} 2 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ -1 \end{bmatrix} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}, \dots \end{aligned}$$

Очевидно, что $x^k = [k(-k+1)]^T$, $x^{k+1} = [(k+1)(-k)]^T$. Тогда $(x^k, x^k) = 2k^2 - 2k + 1$, $(x^k, x^{k+1}) = 2k^2$, и, следовательно,

$$\lambda^{(k+1)} = \frac{2k^2}{2k^2(1 - 1/k + O(k^{-2}))} = 1 + \frac{1}{k} + O\left(\frac{1}{k^2}\right).$$

Сходимость к собственному значению $\lambda = 1$ очень медленна и не похожа на ту, которую мы имели в 2.2. Выясним, с чем это связано. Решая задачу на собственные значения, находим, что рассматриваемая матрица имеет двукратное собственное значение $\lambda = 1$, которому отвечает единственный собственный вектор $[1, -1]^T$. Тем самым, эта матрица не является матрицей простой структуры, как это было в 2.2, а ее жордановой формой является клетка порядка два. Приведенный пример показывает, что степенной метод не отказывается работать и в том случае, когда максимальному по модулю кратному собственному значению отвечает жорданова клетка, но скорость сходимости резко падает от скорости сходимости геометрической прогрессии к скорости сходимости гармонического ряда.

11.5 QR - алгоритм

Кратко опишем один из наиболее употребительных методов при решении полной проблемы собственных значений.

Пусть матрица A разложена в произведение ортогональной Q матрицы и верхней треугольной матрицы R , т.е.

$$A_0 = A = QR = Q_1 R_1.$$

Сделать это можно, например, при помощи метода вращений или метода отражений. Построим матрицу

$$A_1 = R_1 Q_1.$$

Поскольку

$$A_1 = Q_1^{-1}Q_1R_1Q_1 = Q^{-1}AQ, \quad (11.20)$$

то матрицы A и A_1 подобны и, следовательно, имеют одинаковый набор собственных значений.

Далее,

$$A_2 = Q_2R_2, \quad A_3 = R_3Q_3 \quad \text{и т.д.}$$

При некоторых ограничениях на A матрицы A_i по форме сходятся к треугольной матрице, на главной диагонали которой стоят собственные числа. Понимать это нужно так, что поддиагональные элементы стремятся к нулю, диагональные к собственным числам, а наддиагональные могут никуда не стремиться.

Если все $|\lambda_i|$ различны, $a_{ij}^{(k)}$, $i > j$ стремятся к нулю со скоростью геометрической прогрессии со знаменателем (λ_i/λ_j) .

Одна итерация $O(n^3)$. Матрица Хессенберга — почти треугольная матрица (первая поддиагональ отлична от нуля). Одна итерация $O(n^2)$.

Итерации со сдвигом

$$\begin{aligned} A_k - \tau_k I &= Q_k R_k, \\ A_{k+1} &= R_k Q_k + \tau_k I = Q_k^{-1}(A_k - \tau_k I)Q_k + \tau_k I = \\ &= Q_k^{-1}A Q_k. \end{aligned}$$

Итерации очень трудоемкие. Скорость сходимости низкая.

Определение 11.4. Матрица A имеет верхнюю форму Хессенберга, если $a_{ij} = 0$ для любых $i > j + 1$.

Это означает, что матрица имеет почти треугольную форму.

$$\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ & * & * & * & * \\ & & * & * & * \\ & & & * & * \end{bmatrix}$$

Теорема 11.1. Всякая действительная квадратная матрица при помощи ортогонального преобразования подобия может быть приведена к форме Хейсенберга.

Доказательство. Построим требуемое преобразование. Этот алгоритм очень похож на алгоритм QR-разложения матрицы при помощи отражений. Приведение осуществляется за $n - 2$ шага. На первом шаге необходимыми нулями заполнится первый столбец, на втором шаге — второй и т.д.

Выполним первый шаг. Запишем матрицу A в виде

$$A = \begin{bmatrix} a_{11} & c^T \\ b & \hat{A} \end{bmatrix}.$$

Пусть \widehat{U}_1 — матрица отражения, переводящая $(n-1)$ -мерный вектор в вектор $[t \ 0 \ \dots \ 0]^T$, и пусть

$$U_1 = \begin{bmatrix} 1 & 0 \\ 0 & \widehat{U}_1 \end{bmatrix}. \quad (11.21)$$

Тогда

$$U_1 A = \begin{bmatrix} 1 & 0 \\ 0 & \widehat{U}_1 \end{bmatrix} \begin{bmatrix} a_{11} & c^T \\ b & \widehat{A} \end{bmatrix} = \begin{bmatrix} a_{11} & c^T \\ \widehat{U}_1 b & \widehat{U}_1 \widehat{A} \end{bmatrix} = \left[\begin{array}{c|c} a_{11} & c^T \\ \hline t_1 & \\ 0 & \\ \vdots & \\ 0 & \widehat{U}_1 \widehat{A} \end{array} \right].$$

Далее, при вычислении $U_1 A U_1^{-1}$ вспомним, что $U_1^{-1} U_1^T = U_1$, и, следовательно,

$$U_1 A U_1 = \begin{bmatrix} a_{11} & c^T \\ b & \widehat{A} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \widehat{U}_1 \end{bmatrix} = \left[\begin{array}{c|c} a_{11} & c^T \widehat{U}_1 \\ \hline t_1 & \\ 0 & \\ \vdots & \\ 0 & \widehat{U}_1 \widehat{A} \widehat{U}_1 \end{array} \right] = \left[\begin{array}{c|cc} a_{11} & * & \dots & * \\ \hline t_1 & & & \\ 0 & & & \\ \vdots & & & \\ 0 & & & \widehat{A} \widehat{U}_1 \end{array} \right] \quad (11.22)$$

и т.д.

Замечание 11.4. Если бы мы попытались взять такую матрицу отражений U_1 , которая оставляет в первом столбце только один ненулевой элемент, то при умножении $U_1 A$ на U_1 справа нам не удалось бы сохранить структуру первого столбца, т.е. полученные после первого умножения A на U_1 слева нули. Именно благодаря формуле (11.21) для матрицы U_1 операция умножения на U_1 справа не затрагивает нули в первом столбце (11.22).

Замечание 11.5. Для симметричной матрицы A ортогонально подобная ей матрица тоже симметрична и, следовательно, в этом случае матрица Хессенберга будет тридиагональной.

Теорема 11.2. Пусть A_m — невырожденная верхняя хессенбергова матрица и A_{m+1} получена из A_m посредством одной QR-итерации (11.20). Тогда A_{m+1} также имеет верхнюю хессенбергову форму.

Доказательство. Для построения A_{m+1} нам нужно сначала построить разложение $A_m = Q_m R_m$, которое можно переписать в виде $Q_m = A_m R_m^{-1}$. Ранее было показано (упражнение 1.1), что обратная к невырожденной верхней треугольной матрице есть верхняя треугольная матрица. Аналогично доказывается, что произведение верхней треугольной и верхней хессенберговой матрицы в любом порядке есть верхняя хессенбергова матрица. Поэтому Q_m есть верхняя хессенбергова матрица. Но и $A_{m+1} = R_m Q_m$ должна быть верхней хессенберговой. Теорема доказана.

Пусть

$$A_m = \begin{bmatrix} a_{11}^{(m)} & a_{12}^m & \cdots & a_{1,n-1}^{(m)} & a_{1,n}^{(m)} \\ a_{21}^{(m)} & a_{22}^m & \cdots & a_{2,n-1}^{(m)} & a_{2,n}^{(m)} \\ & a_{32}^m & \cdots & a_{3,n-1}^{(m)} & a_{3,n}^{(m)} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ & & & a_{n,n-1}^{(m)} & a_{n,n}^{(m)} \end{bmatrix}$$

Если $|\lambda_i| > |\lambda_{i+1}|$, то

$$a_{i+1,i}^{(m)} = O\left(\left|\frac{\lambda_{i+1}}{\lambda_i}\right|\right)^m.$$

Теорема 11.3 (Теорема Шура). Пусть $A \in \mathbb{C}^n \times \mathbb{C}^n$. Тогда существует унитарная матрица $U \in \mathbb{C}^n \times \mathbb{C}^n$ и верхняя треугольная матрица $T \in \mathbb{C}^n \times \mathbb{C}^n$ такие, что

$$T = U^*AU.$$

$A = UTU^*$ — разложение Шура.

Замечание 11.6. При решении системы (11.18) (или системы с матрицей (11.19)) следует опасаться трудностей, которые могут возникнуть в связи с плохой обусловленностью матрицы A (или $(A - \widehat{\lambda}I)$). С одной стороны, чем ближе к собственному значению, тем быстрее сходятся итерации, с другой стороны, тем хуже обусловленность матрицы, подлежащей обращению.

Глава II

Методы решения нелинейных уравнений

§ 12

Методы решения нелинейных уравнений

Пусть задана непрерывная функция $f(x)$ действительной переменной x , и требуется найти ее нули, т.е. корни уравнения

$$f(x) = 0. \quad (12.1)$$

При такой формулировке задача весьма неопределенна, ибо корней может не быть вовсе, или их может быть бесконечно много. Обычно задача формулируется более конкретно с дополнительными указаниями. Например, отыскание корней на заданном интервале. Поскольку не существует регулярных методов отыскания точных значений корней уравнения (12.1), то речь должна идти об итерационных методах нахождения приближенного решения. (Только если $f(x)$ представляет собой многочлен не выше 4-ой степени, имеются методы представления его нулей в виде радикалов.)

Чтобы воспользоваться тем или иным итерационным методом, нужно иметь начальное приближение к корню. Для этого нужно, по крайней мере, изучить расположение корней и выделить области, где имеется единственный корень. В противном случае мы должны с использованием того или иного итерационного процесса уточнить значения корней или найти их с требуемой точностью.

Способы локализации корней (выделение областей, где имеется единственный корень) многообразны, и указать универсальный метод не представляется возможным. Иногда отрезки локализации известны заранее, а иногда определяются из физических соображений. В простых ситуациях хороший результат может дать графический метод; широко применяют построение таблиц функции $f(x)$ вида $y_i = f(x_i)$, $i = \overline{1, n}$ для обнаружения перемен знака.

12.1 Метод бисекции (метод деления отрезка пополам)

Пусть $f(x) \in C[a, b]$ и $f(a)f(b) < 0$. Последнее означает, что на $[a, b]$ имеется, по крайней мере, один корень уравнения (12.1). (Условие существования решения.) Предположим, что решение единственное, т.е. $x^* \in (a, b)$ — единственный корень уравнения (12.1) на $[a, b]$. Положим $a_0 = a$, $b_0 = b$, найдем середину отрезка $[a_0, b_0]$

$$x_0 = \frac{a_0 + b_0}{2}$$

и примем эту величину за приближенное значение x^* . Так как положение корня x^* на отрезке $[a_0, b_0]$ неизвестно, то можно лишь утверждать, что погрешность этого приближения не превосходит половины длины $[a_0, b_0]$:

$$|x_0 - x^*| \leq \frac{b_0 - a_0}{2}.$$

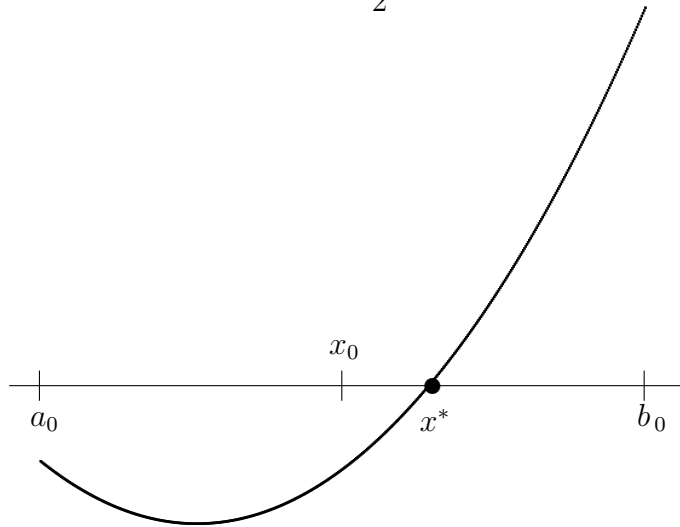


Рис. 1

Вычислим $f(x_0)$. Если $f(x_0) = 0$, то $x^* = x_0$, и вычисления на этом заканчиваются. Если $f(x_0) \neq 0$, то знак $f(x_0)$ совпадает либо со знаком $f(a_0)$, либо со знаком $f(b_0)$. Пусть для определенности $f(a_0) < 0$, $f(b_0) > 0$. Из двух отрезков $[a_0, x_0]$ и $[x_0, b_0]$ выберем тот, на концах которого $f(x)$ принимает значения с противоположными знаками. Обозначим этот отрезок через $[a_1, b_1]$, где

$$a_1 = a_0, \quad b_1 = x_0 \quad \text{при} \quad f(x_0) > 0$$

и

$$a_1 = x_0, \quad b_1 = b_0 \quad \text{при} \quad f(x_0) < 0.$$

Отрезок $[a_1, b_1]$ имеет вдвое меньшую длину, чем $[a_0, b_0]$, $f(a_1)f(b_1) < 0$ и $x^* \in (a_1, b_1)$, причем $|x_0 - x^*| \leq \frac{b_0 - a_0}{2}$. Найдем середину отрезка $[a_1, b_1]$ и т.д. Пусть

$$x_k = \frac{a_k + b_k}{2}, \quad k = 1, 2, \dots,$$

и всегда

$$|x_k - x^*| \leq \frac{b_k - a_k}{2} = \frac{b - a}{2^{k+1}}.$$

Процесс деления отрезка пополам продолжается до тех пор, пока длина нового отрезка $[a_k, b_k]$ не станет меньше 2ε , где ε — требуемая точность в определении приближенного значения корня. Тогда

$$x_k = \tilde{x}, \quad |\tilde{x} - x^*| \leq \varepsilon,$$

т.е. изложенный метод позволяет найти приближенное решение с *гарантированной* точностью. Скорость сходимости метода не ниже скорости сходимости к нулю геометрической прогрессии со знаменателем $1/2$. Каждая итерация уменьшает погрешность не менее, чем в два раза.

Пример 12.1. Для того, чтобы уменьшить первоначальную локализацию в 10^6 раз, нужно сделать 20 итераций, ибо $2^{20} = 1048576$.

12.2 Метод простых итераций

Чтобы применить метод простых итераций для решения нелинейного уравнения (12.1), необходимо преобразовать это уравнение к следующему виду:

$$x = \varphi(x). \quad (12.2)$$

Это можно сделать многими различными способами, некоторые из которых будут изложены позже. Пусть, например,

$$\varphi(x) = x + \tau(x)f(x), \quad (12.3)$$

где $\tau(x)$ — произвольная непрерывная знакоопределенная функция.

Выбирая некоторое начальное приближение x_0 , построим итерационный процесс

$$x_{k+1} = \varphi(x_k), \quad k = 0, 1, \dots \quad (12.4)$$

Итерационный процесс (12.4) называется методом простых итераций.

Теорема 12.1. Пусть x^* — корень уравнения (12.2). Тогда, если $|\varphi'(x)| \leq q < 1$ для $x \in [x^* - \delta, x^* + \delta]$, то при любом начальном приближении $x_0 \in [x^* - \delta, x^* + \delta]$ метод простых итераций сходится со скоростью геометрической прогрессии, знаменателем которой является число q . При этом

$$|x_k - x^*| \leq q^k |x_0 - x^*|.$$

Эта теорема или ей аналогичная была доказана в курсе математического анализа.

Дадим геометрическую иллюстрацию итерационного процесса (12.4). Изобразим на плоскости Oxy прямую $y = x$ и кривую $y = \varphi(x)$. Пусть сначала $0 < \varphi'(x) < 1$.

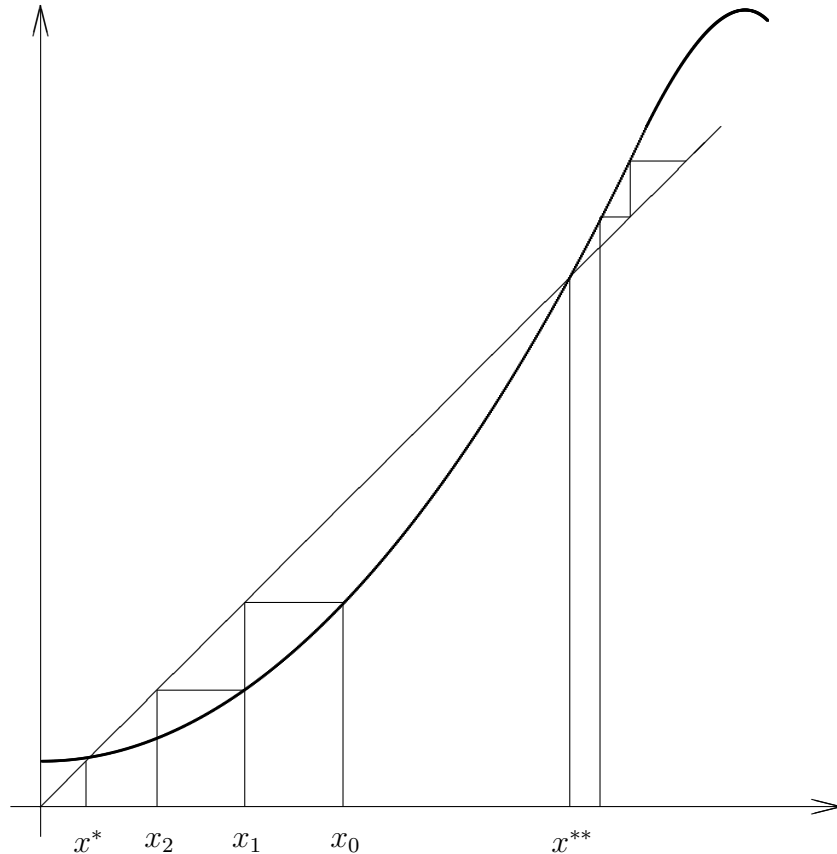


Рис. 2

Из рисунка 2 видно, что при $0 < \varphi'(x) \leq q < 1$ последовательность x_k монотонно сходится к x^* , причем с той стороны, с которой расположено начальное приближение.

Если $\varphi'(x^{**}) > 1$, то итерации не сходятся к x^{**} .

При $-1 < \varphi'(x) < 0$ приближения двусторонние (см. рис. 3). В этом случае по двум последовательным приближениям легко судить о достаточной точности

$$|x_k - x^*| < |x_k - x_{k-1}|.$$

Можно также увидеть, что сходимость тем быстрее, чем меньше $|\varphi'|$. Если $\varphi'(x^{**}) \ll 1$, то процесс сходимости ускоряется по мере приближения к корню.

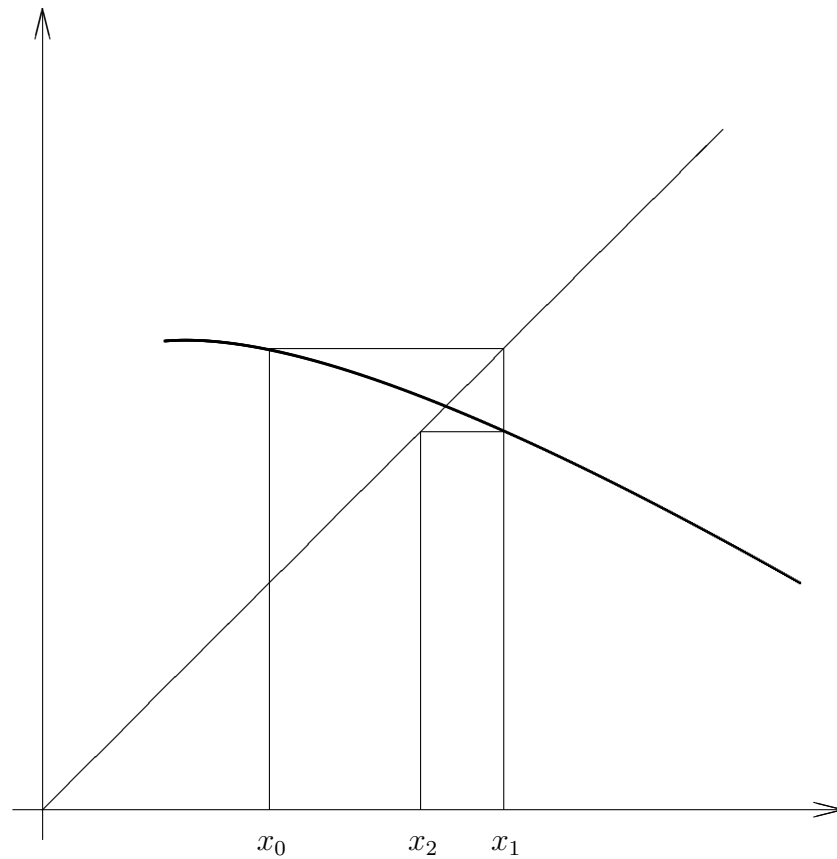


Рис. 3

12.3 Метод Ньютона

Пусть k -е приближение x_k к решению уравнения (12.1) найдено. Разложим функцию $f(x)$ в точке x_k по формуле Тейлора

$$f(x) = f(x_k) + (x - x_k)f'(x_k) + O((x - x_k)^2)$$

или

$$\Delta f = f(x) - f(x_k) = df + O(d^2 f), \quad df \equiv f' dx.$$

Заменим приближенно приращение функции ее дифференциалом

$$\Delta f \approx df$$

или

$$f(x) \approx f(x_k) + (x - x_k)f'(x_k) =: P_1(x).$$

Приравняем теперь функцию $P_1(x)$, являющуюся приближением к $f(x)$, нулю

$$P_1(x) = 0 : \quad f(x_k) + (x - x_k)f'(x_k) = 0$$

и найдем корень полученного уравнения

$$x - x_k = -\frac{f(x_k)}{f'(x_k)} \quad \Rightarrow \quad x = x_k - \frac{f(x_k)}{f'(x_k)}.$$

Этот корень и примем за новое приближение. Итак, алгоритм метода Ньютона следующий:

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, \dots, \quad x_0 \text{ — задано.} \quad (12.5)$$

Геометрическая интерпретация метода Ньютона такова. Как известно из математического анализа, $f'(x_k)$ есть тангенс угла наклона касательной к кривой $y = f(x)$ в точке $x = x_k$. Прямая

$$y = P_1(x) \quad (12.6)$$

имеет тот же наклон, что и касательная к $f(x)$ в точке x_k . Более того, в точке $x = x_k$ значения $P_1(x)$ и $f(x)$ совпадают, и, следовательно, (12.6) есть уравнение касательной к кривой $y = f(x)$ в точке $x = x_k$.

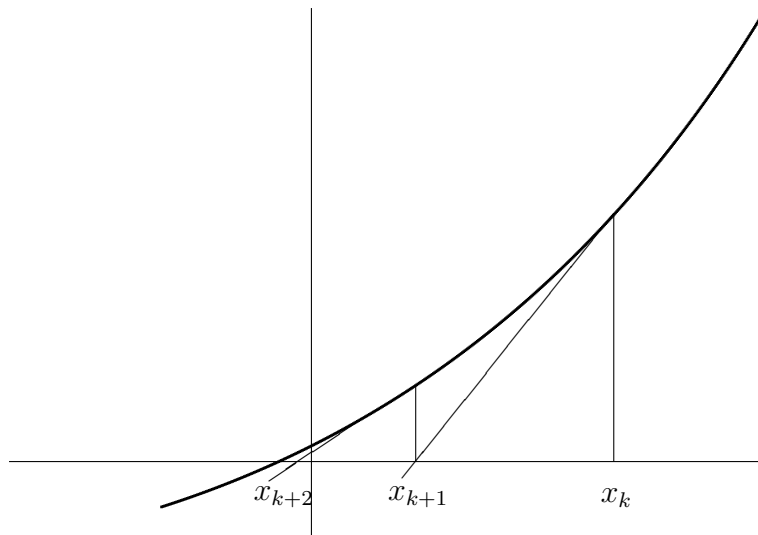


Рис. 4

Тем самым, в методе Ньютона на каждой итерации кривая $y = f(x)$ заменяется касательной в точке x_k , и вместо уравнения $f(x) = 0$ решается уравнение $P_1(x) = 0$.

Связь метода Ньютона с методом простых итераций. При применении метода простых итераций к решению уравнения (12.1) оно сначала преобразовывалось к виду (12.2), где функция $\varphi(x)$ определялась соотношением (12.3), а итерации проводились по формуле (12.4). Сравнивая (12.4) и (12.5), находим, что

$$\varphi(x) = x - \frac{f(x)}{f'(x)}, \quad (12.7)$$

т.е. $\tau(x)$ из (12.3) есть $[-f'(x)]^{-1}$. Как следует из теоремы 12.1, для сходимости метода простых итераций производная функции $\varphi(x)$ в окрестности корня x^* должна быть по модулю меньше единицы. Из (12.7)

$$\varphi'(x) = 1 - \frac{(f'(x))^2 - f(x)f''(x)}{(f'(x))^2} = \frac{f(x)f''(x)}{(f'(x))^2}.$$

Если x^* — простой корень уравнения (12.1), то $f'(x^*) \neq 0$, а $\varphi'(x^*) = 0$, и существует окрестность x^* , где $|\varphi'(x)| < 1$. Поэтому метод Ньютона всегда сходится, если начальное условие выбрано удачно.

Оценка скорости сходимости метода Ньютона. Еще раз разложим $f(x)$ в точке x_k

$$f(x) = f(x_k) + (x - x_k)f'(x_k) + \frac{(x - x_k)^2}{2}f''(\xi_k), \quad \xi_k \in (x, x_k).$$

Полагая здесь $x = x^*$, получим

$$0 = f(x_k) + (x^* - x_k)f'(x_k) + \frac{(x^* - x_k)^2}{2}f''(\xi_k^*), \quad \xi_k^* \in (x^*, x_k).$$

В силу (12.5)

$$0 = (x_{k+1} - x_k)f'(x_k) + f(x_k).$$

Вычитая это соотношение из предыдущего, получим

$$0 = (x^* - x_{k+1})f'(x_k) + \frac{1}{2}(x^* - x_k)f''(\xi_k^*).$$

Отсюда

$$(x_{k+1} - x^*) = \frac{1}{2} \frac{f''(\xi_k^*)}{f'(x_k)} (x_k - x^*)^2. \quad (12.8)$$

Будем предполагать, что

$$|f'(x)| \geq m_1, \quad |f''(x)| \leq M_2. \quad (12.9)$$

Тогда

$$|x_{k+1} - x^*| \leq \frac{M_2}{2m_1} |x_k - x^*|^2.$$

Домножая левую и правую части этого неравенства на $M_2/(2m_1)$, получим, что

$$\beta_{k+1} = \frac{M_2}{2m_1}|x_{k+1} - x^*| \leq \left[\frac{M_2}{2m_1}(x_k - x^*) \right]^2 = \beta_k^2,$$

т.е.

$$\beta_{k+1} \leq \beta_k^2, \quad (12.10)$$

где

$$\beta_k = \frac{M_2}{2m_1}|x_k - x^*|. \quad (12.11)$$

Из (12.10) имеем

$$\beta_1 \leq \beta_0^2, \quad \beta_2 \leq \beta_1^2 \leq \beta_0^4 = \beta_0^{2^2}, \quad \beta_3 \leq \beta_2^2 \leq \beta_0^{2^3}$$

и вообще

$$\beta_k^2 \leq \beta_0^{2^k}.$$

Принимая во внимание (12.11), находим, что

$$|x_k - x^*| \leq \frac{2m_1}{M_2} \left[\frac{M_2}{2m_1}|x_0 - x^*| \right]^{2^k}. \quad (12.12)$$

Для сходимости нужно, чтобы

$$\frac{M_2}{2m_1}|x_0 - x^*| \leq q < 1. \quad (12.13)$$

Итак, доказана

Теорема 12.2. Пусть $f(x) \in C^2[x^* - \delta, x^* + \delta]$, где x^* — простой корень уравнения (12.1), и при $x \in [x^* - \delta, x^* + \delta]$ справедливы оценки (12.9). Тогда, если начальное приближение $x_0 \in [x^* - \delta, x^* + \delta]$ таково, что справедливо (12.13), то метод Ньютона (12.5) сходится с квадратичной скоростью, и справедлива оценка (12.12).

Пример 12.2. Пусть требуется найти корень степени p из числа $a > 0$. Тогда

$$\begin{aligned} f(x) &:= x^p - a \\ x_{k+1} &= x_k - \frac{x_k^p - a}{px_k^{p-1}} = \frac{p-1}{p}x_k + \frac{a}{px_k^{p-1}}. \end{aligned} \quad (12.14)$$

При $p = 2$ и $a = 3$

$$x_{k+1} = \frac{x_k}{2} + \frac{3}{2x_k}.$$

Пусть $x_0 = 2$. Тогда

$$x_1 = \frac{7}{4} = 1.75, \quad x_2 = \frac{97}{56} \approx 1.7321, \quad x_3 \approx 1.7320508,$$

а $x^* = 1.7320508\dots$. Отсюда следует, что если в x_0 один верный знак, то в x_1 — два, в x_2 — четыре и т.д. Грубо говоря, число верных знаков после каждой итерации удваивается.

Замечание 12.1. Сходимость метода Ньютона установлена при условии, что корень x^* является простым. Ну, а что будет, если корень окажется кратным? Чтобы ответить на этот вопрос, исследуем $ff''/(f')^2$. Если корень x^* имеет кратность $p > 1$, то

$$\begin{aligned} f(x) &= a(x - x^*)^p + O((x - x^*)^{p+1}) \\ f'(x) &= ap(x - x^*)^{p-1} + O((x - x^*)^p) \\ f''(x) &= ap(p-1)(x - x^*)^{p-2} + O((x - x^*)^{p-1}). \end{aligned}$$

Отсюда

$$\begin{aligned} \varphi'(x) &= \frac{f(x)f''(x)}{[f'(x)]^2} = \\ &= \frac{a^2p(p-1)(x - x^*)^{2p-2} + O((x - x^*)^{2p-1})}{p^2a^2(x - x^*)^{2p-2} + O((x - x^*)^{2p-1})} = \frac{p-1}{p} + O(x - x^*) \end{aligned} \quad (12.15)$$

и в малой окрестности x^* $|\varphi'(x)| < 1$. Тем самым, метод Ньютона будет сходиться и к кратному корню, но эта сходимость не будет квадратичной; она будет скоростью сходимости геометрической прогрессии со знаменателем $q = (p-1)/p < 1$.

Возникает вопрос, а нельзя ли увеличить скорость сходимости к кратному корню? Ответ на этот вопрос положительный. Это можно сделать путем следующего обобщения метода Ньютона. Итерации будем вести по формуле

$$x_{k+1} = x_k - \tau \frac{f(x_k)}{f'(x_k)},$$

где значение параметра τ определяется кратностью искомого корня. Найдем это значение. Имеем

$$\varphi(x) = x - \tau \frac{f(x)}{f'(x)}, \quad \varphi'(x) = 1 - \tau \frac{f'^2 - ff''}{f'^2} = (1 - \tau) + \tau \frac{ff''}{f'^2}.$$

Отсюда с учетом (12.15)

$$\varphi'(x) = (1 - \tau) + \tau \left[\frac{p-1}{p} + O(x - x^*) \right] = 1 - \tau + \tau \frac{p-1}{p} + O(x - x^*).$$

Выберем τ из условия, что $\varphi'(x^*) = 0$. Тогда $\tau = p$, и обобщенный метод Ньютона

$$x_{k+1} = x_k - p \frac{f(x_k)}{f'(x_k)},$$

где p — кратность искомого корня, обладает скоростью сходимости метода Ньютона к простому корню.

Пример 12.3. Пусть $f(x) = x^2$. Здесь $x^* = 0$ есть двукратный корень. Метод Ньютона приводит к соотношению

$$x_{k+1} = x_k - \frac{x_k^2}{2x_k} = \frac{1}{2}x_k.$$

Эти соотношения представляют собой геометрическую прогрессию со знаменателем $q = 1/2$.

Обобщенный метод для этого примера сразу дает точное решение.

Пример 12.4. Пусть $f(x) = x^2(x + 1)$ и $x^* = 0$ — двукратный корень. Обобщенный метод Ньютона принимает вид

$$x_{k+1} = \frac{x_k^2}{2 + 3x_k}.$$

Если $x_0 = 1$, то $x_1 = \frac{1}{5} = 0.2$, $x_2 = \frac{1}{65} = 0.015$, $x_3 = 0.0000115$ и т.д.

§ 13

Метод секущих

Основным достоинством метода Ньютона, которое делает его очень привлекательным, является высокая скорость сходимости. К недостаткам следует отнести необходимость вычисления на каждом шаге итераций производной. Вторым недостатком — сильная зависимость результативности метода от начального приближения: если начальное приближение оказалось неудачным (недостаточно близким к искомому решению), метод просто расходится. Первый недостаток в какой-то мере может быть преодолен путем замены производной разностным отношением. Именно, заменим в (13.5) производную $f'(x_k)$ на разностное отношение

$$\frac{\Delta f_k}{\Delta x_k} = \frac{f_k - f_{k-1}}{x_k - x_{k-1}},$$

где $f_k = f(x_k)$. В результате будем иметь

$$x_{k+1} = x_k - f_k \frac{x_k - x_{k-1}}{f_k - f_{k-1}}, \quad k = 1, 2, \dots \quad (13.1)$$

Отметим, что этот метод двухшаговый: чтобы найти x_{k+1} , нужно знать x_k и x_{k-1} . В частности, для того, чтобы начать итерации, также требуются значения начальных приближений x_0 и x_1 .

Обратимся к геометрической интерпретации метода (13.1). Пусть

$$L_1(x) = f_{k-1} \frac{x - x_k}{x_{k-1} - x_k} + f_k \frac{x - x_{k-1}}{x_k - x_{k-1}} \quad (13.2)$$

— интерполяционный многочлен Лагранжа первой степени, построенный по значениям функции $f(x)$ в узлах x_{k-1} и x_k . Рассмотрим прямую

$$y = L_1(x)$$

и найдем ее нуль, т.е. решение уравнения $L_1(x) = 0$. Принимая во внимание (13.2), находим, что

$$f_{k-1}(x - x_k) = f_k(x - x_{k-1}) \equiv f_k(x - x_k) + f_k(x_k - x_{k-1}).$$

Отсюда

$$x = x_k - f_k \frac{x_k - x_{k-1}}{f_k - f_{k-1}},$$

что совпадает с x_{k+1} из (13.1). Отсюда следует, что если в методе Ньютона кривая $y = f(x)$ всякий раз заменяется касательной в точке x_k , то в методе (13.1) кривая $y = f(x)$ заменяется секущей, пересекающей $y = f(x)$ при $x = x_{k-1}$ и $x = x_k$. Эта геометрическая интерпретация метода (13.1) и дает ему название — метод секущих.

Обратимся к оценке скорости сходимости метода секущих. Сначала заметим, что в силу (13.2) и с учетом формулы конечных приращений

$$\begin{aligned} L_1(x_{k+1}) - L_1(x^*) &= \\ &= f_{k-1} \frac{x_{k+1} - x_k}{x_{k-1} - x_k} + f_k \frac{x_{k+1} - x_{k-1}}{x_k - x_{k-1}} - f_{k-1} \frac{x^* - x_k}{x_{k-1} - x_k} - f_k \frac{x^* - x_{k-1}}{x_k - x_{k-1}} = \\ &= f_{k-1} \frac{x_{k+1} - x^*}{x_{k-1} - x_k} + f_k \frac{x_{k+1} - x^*}{x_k - x_{k-1}} = \frac{f_k - f_{k-1}}{x_k - x_{k-1}} (x_{k+1} - x^*) = \\ &= f'(\xi_k)(x_{k+1} - x^*), \quad \xi_k \in (x_{k-1}, x_k) \end{aligned}$$

и поэтому

$$x_{k+1} - x^* = \frac{L_1(x_{k+1}) - L_1(x^*)}{f'(\xi_k)} = -\frac{L_1(x^*)}{f'(\xi_k)}, \quad (13.3)$$

ибо по построению $L_1(x_{k+1}) = 0$. При вычислении $L_1(x^*)$ воспользуемся формулой для погрешности интерполяции

$$f(x) - L_1(x) = \frac{1}{2} f''(\eta_k)(x - x_{k-1})(x - x_k), \quad \eta_k \in (x, x_{k-1}, x_k),$$

где (x, x_{k-1}, x_k) — открытый интервал, концами которого являются крайние из указанных трех точек. Полагая здесь $x = x^*$ и принимая во внимание, что $f(x^*) = 0$, находим искомое выражение, подставляя которое в (13.3), с учетом вышесказанного будем иметь

$$x_{k+1} - x^* = \frac{1}{2} \frac{f''(\eta_k)}{f'(\xi_k)} (x_k - x^*)(x_{k-1} - x^*). \quad (13.4)$$

Как и при изучении метода Ньютона, будем предполагать, что

$$|f'(x)| \geq m_1, \quad |f''(x)| \leq M_2. \quad (13.5)$$

Тогда из (13.4) будем иметь оценку

$$|x_{k+1} - x^*| \leq \frac{M_2}{2m_1} |x_k - x^*| |x_{k-1} - x^*|.$$

Домножая теперь обе части этого неравенства на $M_2/(2m_1)$ и обозначая

$$\frac{M_2}{2m_1}|x_k - x^*| = \Delta_k, \quad (13.6)$$

найдем, что

$$\Delta_{k+1} \leq \Delta_k \Delta_{k-1}. \quad (13.7)$$

Отметим, что в методе Ньютона $\Delta_{k+1} \leq \Delta_k^2$, и это позволило нам говорить о квадратичной сходимости метода. Для метода секущих аналогичное неравенство

$$\Delta_{k+1} \leq \Delta_k^\nu \quad (13.8)$$

доказывается существенно более сложно. Мы не будем его доказывать, а получим оценку типа (13.8) для мажоранты погрешности.

Заметим, что, если

$$z_{k+1} = z_k z_{k-1}, \quad z_0 \geq \Delta_0, \quad z_1 \geq \Delta_1, \quad (13.9)$$

то

$$\Delta_{k+1} \leq z_{k+1}. \quad (13.10)$$

Попытаемся свести нелинейное разностное уравнение (13.9) к виду

$$z_{k+1} = z_k^\nu. \quad (13.11)$$

Если это так, то $z_k = z_{k-1}^\nu$, и, следовательно,

$$z_{k-1} = z_k^{1/\nu}.$$

Подставляя это соотношение и (13.11) в (13.9), получим

$$z_k^\nu = z_k^{1+1/\nu}.$$

Приравнявая степени z_k в левой и правой частях, находим, что

$$\nu = 1 + 1/\nu, \quad \text{т.е.} \quad \nu^2 - \nu - 1 = 0. \quad (13.12)$$

Отсюда вытекает, что

$$\nu_{1,2} = \frac{1 \pm \sqrt{5}}{2}.$$

Корню $\nu = (1 - \sqrt{5})/2$ отвечает неубывающее решение уравнения (13.11), которое не может описывать сходящийся итерационный процесс. Поэтому следует взять

$$\nu = \frac{1 + \sqrt{5}}{2} \approx 1.6180339. \quad (13.13)$$

Итак, вместо (13.8) имеем (13.10), (13.11), (13.13), что позволяет говорить о сходимости метода секущих со скоростью (13.13). Эта скорость меньше, чем у метода Ньютона.

Пример 13.1. Пусть требуется найти корень функции $f(x) = x^2 - a$. Метод секущих применительно к этой функции принимает вид

$$x_{k+1} = x_k - \frac{x_k^2 - a}{x_k^2 - x_{k-1}^2} = \frac{x_k x_{k-1} + a}{x_k + x_{k-1}}.$$

Если $a = 3$ и положить $x_0 = 3$, а $x_1 = 2$, то

$$x_2 = 1.8, \quad x_3 = \frac{33}{19} \approx 1.7368, \quad x_4 \approx 1.7321428$$

при $x^* = 1.732051 \dots$

Замечание 13.1. Нелинейное разностное уравнение (13.11) имеет очевидное решение

$$z_k = z_0^{\nu^k}, \quad (13.14)$$

для которого, в частности, $z_1 = z_0^{\nu}$. Но в итерационном методе (13.1) используются два начальных условия, и поэтому величина z_1 не должна зависеть от z_0 . Мы вынуждены констатировать, что найденное решение (13.14) не совсем правильно описывает этот процесс. То, что мы не получили решения уравнения (13.9), удовлетворяющего обоим начальным условиям, не должно вызывать удивления: мы ведь нашли решение, общее для (13.9) и (13.11), а интересующее нас решение может (13.11) и не удовлетворять.

Вернемся к задаче (13.9) и найдем ее решение. Логарифмируя уравнение (13.9), будем иметь

$$\ln z_{k+1} = \ln z_k + \ln z_{k-1}.$$

Обозначая

$$\ln z_k = y_k, \quad (13.15)$$

для y_k получим линейное разностное уравнение с постоянными коэффициентами

$$y_{k+1} = y_k + y_{k-1}.$$

Его характеристическое уравнение есть

$$q^2 - q - 1 = 0$$

(сравни с (13.12)) с корнями

$$q_1 = \frac{1 + \sqrt{5}}{2}, \quad q_2 = \frac{1 - \sqrt{5}}{2}, \quad q_1 + q_2 = 1, \quad \frac{q_2}{q_1} = \frac{-3 + \sqrt{5}}{2} \approx -0.38.$$

Поэтому

$$y_k = c_1 q_1^k + c_2 q_2^k. \quad (13.16)$$

Удовлетворяя начальным условиям (13.9), будем иметь

$$\begin{aligned} c_1 + c_2 &= y_0 = \ln z_0, \\ q_1 c_1 + q_2 c_2 &= y_1 = \ln z_1. \end{aligned}$$

Отсюда находим, что

$$\begin{aligned} c_1 &= \frac{y_1 - y_0 q_2}{\sqrt{5}} = \frac{y_1 - y_0 + y_0 q_1}{\sqrt{5}}, \\ c_2 &= -\frac{y_1 - y_0 q_1}{\sqrt{5}} = -\frac{y_1 - y_0 + y_0 q_2}{\sqrt{5}} \end{aligned}$$

и, следовательно,

$$\begin{aligned} y_k &= \frac{y_1 + (q_1 - 1)y_0}{\sqrt{5}} q_1^k - \frac{y_1 + (q_2 - 1)y_0}{\sqrt{5}} q_2^k = \\ &= \frac{\ln(z_1 z_0^{q_1 - 1})}{\sqrt{5}} q_1^k - \frac{\ln(z_1 z_0^{q_2 - 1})}{\sqrt{5}} q_2^k, \end{aligned}$$

а с учетом (13.15)

$$\begin{aligned} z_k &= \exp \left\{ \ln(z_1 z_0^{q_1 - 1}) \frac{q_1^k}{\sqrt{5}} - \ln(z_1 z_0^{q_2 - 1}) \frac{q_2^k}{\sqrt{5}} \right\} = \\ &= (z_1 z_0^{q_1 - 1})^{q_1^k / \sqrt{5}} / (z_1 z_0^{q_2 - 1})^{q_2^k / \sqrt{5}}. \end{aligned}$$

Отсюда, в частности, следует, что, если

$$z_1 = z_0^{q_1},$$

то

$$z_1 z_0^{q_2 - 1} = z_0^{q_1 + q_2 - 1} = 1, \quad z_1 z_0^{q_1 - 1} = z_0^{2q_1 - 1} = z_0^{\sqrt{5}}$$

и, следовательно,

$$z_k = z_0^{q_1^k},$$

что совпадает с (13.14), ибо $q_1 = \nu$.

В более же реалистическом случае, когда $z_1 = z_0$,

$$z_k = z_0^{(q_1^{k+1} - q_2^{k+1}) / \sqrt{5}} = z_0^{q_1^{k+1} (1 - (q_2/q_1)^{k+1}) / \sqrt{5}}.$$

Это соотношение дает представление погрешности на k -ой итерации через погрешности z_0 и $z_1 = z_0$.

Небезынтересен вопрос о скорости убывания погрешности на двух соседних итерациях. Пусть

$$z_{k+1} = z_k^{\nu_k}.$$

Отсюда

$$\nu_k = (\ln z_{k+1}) / (\ln z_k) = y_{k+1} / y_k,$$

а с учетом (13.16)

$$\begin{aligned} \nu_k &= \frac{y_{k+1}}{y_k} = \frac{c_1 q_1^{k+1} + c_2 q_2^{k+1}}{c_1 q_1^k + c_2 q_2^k} = \\ &= q_1 \frac{1 + \frac{c_2}{c_1} (q_2/q_1)^{k+1}}{1 + \frac{c_2}{c_1} (q_2/q_1)^k} = q_1 \left[1 - \frac{c_2}{c_1} \left(\frac{q_2}{q_1} \right)^k + O \left(\frac{q_2}{q_1} \right)^{k+1} \right] = \\ &= q_1 + O \left((q_2/q_1)^k \right). \end{aligned}$$

Отсюда следует, что при любых начальных приближениях убывание погрешности в малой окрестности корня происходит со скоростью $\sim q_1$.

Замечание 13.2. Мы уже отмечали, что скорость сходимости метода секущих ниже, чем метода Ньютона. И, тем не менее, метод секущих может оказаться более предпочтительным по сравнению с методом Ньютона. Для реализации каждой итерации в методе Ньютона нужно вычислять значение функции и ее производной в новой точке. В методе секущих на каждой итерации нужно знать только одно новое значение функции. Если эти операции трудоемкие, то две операции по методу секущих могут быть сравнимы по трудоемкости с одной итерацией по методу Ньютона, а это приводит к большему уменьшению начальной погрешности.

Глава III

Численное дифференцирование

§ 14

Численное дифференцирование

14.1 Введение

Численное дифференцирование применяется, если функция задана таблицей или если ее трудно продифференцировать аналитически. Допустим, что в окрестности некоторой точки x у функции $f(x)$ существует производная. По определению

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}.$$

Если отказаться от предельного перехода, то можно положить

$$f'(x) \approx \frac{f(x + \Delta x) - f(x)}{\Delta x}. \quad (14.1)$$

Это и есть простейшая формула численного дифференцирования. Оценим ее погрешность в предположении, что значения функции $f(x)$ вычисляются точно, и она дважды непрерывно дифференцируема. Используя формулу Тейлора, находим, что

$$\begin{aligned} \frac{f(x + \Delta x) - f(x)}{\Delta x} &= \frac{f(x) + \Delta x f'(x) + \frac{(\Delta x)^2}{2} f''(\xi) - f(x)}{\Delta x} = \\ &= f'(x) + \frac{\Delta x}{2} f''(\xi), \quad \xi \in (x, x + \Delta x). \end{aligned} \quad (14.2)$$

Отсюда заключаем, что формула (14.1) для функции $f(x) \in C^2$ имеет погрешность первого порядка малости относительно Δx .

Пусть $x_i = x_0 + ih$, где $i \in \mathbb{Z}$, а $h > 0$ — шаг сетки. Тогда, полагая в (14.2) $x = x_i$, а $\Delta x = h$, получим

$$\frac{f(x_{i+1}) - f(x_i)}{h} = f'(x_i) + \frac{h}{2} f''(\xi_i). \quad (14.3)$$

Если же в (14.2) положить $\Delta x = -h$ и снова $x = x_i$, то

$$\frac{f(x_{i-1}) - f(x_i)}{-h} = f'(x_i) - \frac{h}{2} f''(\tilde{\xi}_i). \quad (14.4)$$

Из (14.3), (14.4) вытекает, что и приближенная формула

$$f'(x_i) \approx \frac{f(x_{i+1}) - f(x_i)}{h} \quad (14.5)$$

и приближенная формула

$$f'(x_i) \approx \frac{f(x_i) - f(x_{i-1}))}{h} \quad (14.6)$$

являются формулами первого порядка точности, однако их погрешности, вообще говоря, имеют разные знаки. Поэтому есть надежда, что у полусуммы правых частей (14.5), (14.6) погрешность будет иметь бóльший порядок малости относительно h (при бóльшей гладкости). В самом деле, используя формулу Тейлора, находим, что

$$\begin{aligned} & \frac{1}{2} \left[\frac{f(x_{i+1}) - f(x_i)}{h} + \frac{f(x_i) - f(x_{i-1}))}{h} \right] = \frac{f(x_{i+1}) - f(x_{i-1}))}{2h} = \\ & = \left[f_i + hf'_i + \frac{h^2}{2}f''_i + \frac{h^3}{6}f'''(\bar{\xi}_i) - \left(f_i - hf'_i + \frac{h^2}{2}f''_i - \frac{h^3}{6}f'''(\bar{\xi}_i) \right) \right] / 2h \\ & = f'_i + \frac{h^2}{6} \frac{f'''(\bar{\xi}_i) + f'''(\bar{\xi}_i)}{2} = f'_i + \frac{h^2}{6} f'''(\xi_i), \quad \xi_i \in (\bar{\xi}_i, \bar{\xi}_i) \subset (x_{i-1}, x_{i+1}). \end{aligned} \quad (14.7)$$

Отсюда заключаем, что для $f(x) \in C^3$ формула

$$f'(x_i) \approx \frac{f(x_{i+1}) - f(x_{i-1}))}{2h} \quad (14.8)$$

имеет погрешность $O(h^2)$.

Теперь вычтем из (14.3) соотношение (14.4)

$$\frac{f(x_{i+1}) - 2f(x_i) + f(x_{i-1}))}{h} = h \frac{f''(\xi_i) + f''(\tilde{\xi}_i)}{2} = f''(\tilde{\xi}_i)h.$$

Следовательно,

$$\frac{f_{i+1} - 2f_i + f_{i-1}}{h^2} = f''(\tilde{\xi}_i), \quad \tilde{\xi}_i \in (x_{i-1}, x_{i+1}), \quad (14.9)$$

т.е. левая часть этого соотношения аппроксимирует вторую производную функции $f(x)$. Исследуем погрешность этой аппроксимации в точке x_i :

$$\begin{aligned} & \frac{f_{i+1} - 2f_i + f_{i-1}}{h^2} = \frac{1}{h^2} \left[f_i + hf'_i + \frac{h^2}{2!}f''_i + \frac{h^3}{3!}f'''_i + \frac{h^4}{4!}f^{IV}(\bar{\xi}_i) - \right. \\ & \left. - 2f_i + f_i - hf'_i + \frac{h^2}{2!}f''_i - \frac{h^3}{3!}f'''_i + \frac{h^4}{4!}f^{IV}(\bar{\xi}_i) \right] = \\ & = f''_i + \frac{h^2}{12} \frac{f^{IV}(\bar{\xi}_i) + f^{IV}(\bar{\xi}_i)}{2} = f''(x_i) + \frac{h^2}{12} f^{IV}(\xi_i). \end{aligned} \quad (14.10)$$

Отсюда следует, что левая часть соотношения (14.10) аппроксимирует вторую производную функции $f(x) \in C^4$ с погрешностью $O(h^2)$.

Замечание 14.1. Мы уже трижды (в (14.7), (14.10) и в формуле после соотношения (14.8)) воспользовались утверждением о том, что для непрерывной функции $0.5(f(x) + f(y)) = f(z)$, где $z \in (x, y)$. Докажем это утверждение в более общем виде. Пусть $f(x) \in C[a, b]$, $m = \min_{[a, b]} f(x)$, $M = \max_{[a, b]} f(x)$, $x_1, x_2 \in [a, b]$, $\alpha > 0$, $\beta > 0$. Тогда

$$\eta = \frac{\alpha f(x_1) + \beta f(x_2)}{\alpha + \beta} = f(x_3), \quad x_3 \in [a, b].$$

В самом деле,

$$m \leq \frac{\alpha f(x_1) + \beta f(x_2)}{\alpha + \beta} = \eta \leq M.$$

и по теореме о промежуточных значениях $\eta = f(x_3)$.

Введем следующие обозначения

$$f_{x,i} = \frac{f_{i+1} - f_i}{h}, \quad f_{\bar{x},i} = \frac{f_i - f_{i-1}}{h}, \quad f_{\circ,x,i} = \frac{1}{2}(f_{x,i} + f_{\bar{x},i}).$$

Тогда

$$f_{\bar{x}x,i} = \frac{f_{i+1} - 2f_i + f_{i-1}}{h^2}.$$

Отсюда и из (14.3), (14.4), (14.7), (14.10) находим, что

$$\begin{aligned} f_{x,i} &= f'_i + O(h), \\ f_{\bar{x},i} &= f'_i + O(h), \\ f_{\circ,x,i} &= f'_i + O(h^2), \\ f_{\bar{x}x,i} &= f''_i + O(h^2). \end{aligned} \tag{14.11}$$

14.2 Метод неопределенных коэффициентов

Рассмотренные простейшие формулы численного дифференцирования были построены из неких эвристических соображений. Существуют и регулярные способы построения формул численного дифференцирования. Один из них — метод неопределенных коэффициентов.

Будем искать формулу численного нахождения k -ой производной в следующем виде

$$f^{(k)}(x) \approx \sum_{j=0}^n c_j f(x_j), \quad k \leq n \tag{14.12}$$

и выберем c_j из тех условий, чтобы формула была точна на многочленах некоторой степени. Рассмотрим

Пример 14.1. Пусть

$$f'(h) \approx c_0 f(0) + c_1 f(h) + c_2 f(2h). \quad (14.13)$$

Потребуем, чтобы формула была точна на линейных функциях. Подставляя в (14.13) $f(x) \equiv 1$ и $f(x) \equiv x$ и требуя выполнения точного равенства, будем иметь

$$\begin{aligned} 0 &= c_0 + c_1 + c_2, \\ 1 &= c_1 h + c_2 2h. \end{aligned}$$

Принимая c_0 за параметр, для c_1 и c_2 получим систему

$$\begin{aligned} c_1 + c_2 &= -c_0, \\ hc_1 + 2hc_2 &= 1. \end{aligned} \quad (14.14)$$

Определитель этой системы равен h и поэтому

$$\begin{aligned} c_1 &= \begin{vmatrix} -c_0 & 1 \\ 1 & 2h \end{vmatrix} / h = (-2hc_0 - 1)/h, \\ c_2 &= \begin{vmatrix} 1 & -c_0 \\ h & 1 \end{vmatrix} / h = (1 + c_0 h)/h. \end{aligned} \quad (14.15)$$

Итак, мы построили однопараметрическое семейство трехточечных формул численного нахождения первой производной

$$f'(h) \approx c_0 f(0) - \frac{1 + 2c_0 h}{h} f(h) + \frac{1 + c_0 h}{h} f(2h).$$

При $c_0 = 0$ имеем

$$f'(h) \approx \frac{f(2h) - f(h)}{h} = f_x(h).$$

При $c_0 = -1/h$

$$f'(h) \approx \frac{f(h) - f(0)}{h} = f_{\bar{x}}(h).$$

Это уже известные нам формулы.

Потребуем теперь, чтобы (14.13) была точна на многочленах второй степени. Тогда к уравнениям (14.14) добавится еще одно уравнение

$$2h = c_1 h^2 + c_2 4h^2.$$

Подставляя сюда c_1 и c_2 из (14.15), получим

$$-(1 + 2c_0 h)h + 4(1 + c_0 h)h = 2h,$$

откуда находим $c_0 = -1/2h$. Подставляя это значение в (14.15), будем иметь $c_1 = 0$, $c_2 = 1/2h$, и поэтому

$$f'(h) \approx \frac{f(2h) - f(0)}{2h} \equiv f_{\bar{x}}(h).$$

И эта формула нам уже известна.

Построим теперь новую формулу.

Пример 14.2. Пусть теперь (ср. с (14.13))

$$f'(2h) \approx c_0 f(0) + c_1 f(h) + c_2 f(2h). \quad (14.16)$$

Будем требовать, чтобы формула (14.16) была точна на многочленах второй степени. Подставляя в (14.16) последовательно $f(x) \equiv 1$, $f(x) \equiv x$ и $f(x) \equiv x^2$ и требуя выполнения точного равенства, получим систему

$$\begin{aligned} c_0 + c_1 + c_2 &= 0, \\ hc_1 + 2hc_2 &= 1, \\ h^2c_1 + 4h^2c_2 &= 4h. \end{aligned} \quad (14.17)$$

Первые два уравнения (14.17) совпадают с (14.14). Поэтому c_1 и c_2 выражаются через c_0 при помощи (14.15). Подставляя (14.15) в последнее уравнение (14.17), находим, что

$$-(2hc_0 + 1)h + 4h(1 + hc_0) = 4h$$

и, следовательно, $c_0 = 1/2h$, а с учетом (14.15)

$$c_1 = -2/h, \quad c_2 = 3/(2h).$$

Тем самым,

$$f'(2h) \approx \frac{f_0 - 4f_1 + 3f_2}{2h} = \frac{f_2 - f_1}{h} + \frac{f_0 - 2f_1 + f_2}{2h} = f_{\bar{x},2} + \frac{h}{2} f_{\bar{x}\bar{x},2}. \quad (14.18)$$

Эта новая формула для вычисления первой производной.

Упражнение 14.1. Показать, что для $f(x) \in C^3$ формула (14.18) имеет погрешность $O(h^2)$.

14.3 Использование интерполяционных формул

Наиболее универсальный способ построения формул численного дифференцирования основан на использовании интерполяционных формул.

Как известно,

$$f(x) = L_n(x) + R_n(x), \quad L_n(x) \equiv \sum_{i=0}^n f_i \prod_{k \neq i} \frac{x - x_k}{x_i - x_k},$$

где $L_n(x)$ — интерполяционный многочлен Лагранжа степени n , а $R_n(x)$ — остаточный член. Полагая

$$f^{(m)}(x) \approx L_n^{(m)}(x), \quad 0 \leq m \leq n, \quad (14.19)$$

получим формулу численного дифференцирования.

Пример 14.3. Пусть $n = 1$, $x_1 = x_0 + h$. Тогда

$$L_1(x) = f_0 \frac{x - x_1}{-h} + f_1 \frac{x - x_0}{h}, \quad L_1'(x) = \frac{f_0}{-h} + \frac{f_1}{h} = \frac{f_1 - f_0}{h}.$$

Полагая здесь $x = 0$, получим формулу (14.5), полагая $x = h$ — формулу (14.6).

Пример 14.4. Пусть теперь $n = 2$, $x_0 = 0$, $x_1 = h$, $x_2 = 2h$. Тогда

$$L_2(x) = f_0 \frac{x - x_1}{x_0 - x_1} \frac{x - x_2}{x_0 - x_2} + f_1 \frac{x - x_0}{x_1 - x_0} \frac{x - x_2}{x_1 - x_2} + f_2 \frac{x - x_0}{x_2 - x_0} \frac{x - x_1}{x_2 - x_1},$$

$$L_2'(x) = f_0 \frac{2x - (x_2 + x_1)}{(-h)(-2h)} + f_1 \frac{2x - (x_2 + x_0)}{h(-h)} + f_2 \frac{2x - (x_0 + x_1)}{2h \cdot h}.$$

Отсюда

$$L_2'(x_2) = L_2'(2h) = \frac{f_0 - 4f_1 + 3f_2}{2h} \quad (\text{ср. с (14.18)})$$

$$L_2'(x_1) = L_2'(h) = \frac{f_2 - f_0}{2h} \quad (\text{ср. с (14.8)})$$

$$L_2'(x_0) = L_2'(0) = \frac{-3f_0 + 4f_1 - f_2}{2h}.$$

Это новая формула; ее погрешность на функциях из C^3 есть $O(h^2)$.

Пример 14.5. Пусть теперь $n = 2$, а $x_0 = 0$, $x_1 = h_1$, $x_2 = x_1 + h_2$, $m = 2$. Легко проверить, что

$$L_2''(x) = \frac{2}{-h_1(-h_1 - h_2)} f_0 + \frac{2}{h_1(-h_2)} f_1 + \frac{2}{(h_1 + h_2)h_2} f_2.$$

Если $h_1 = h_2 = h$, то

$$L_2''(x) = \frac{f_0 - 2f_1 + f_2}{h^2} = f_{\bar{x}x}(h).$$

В противном случае

$$L_2''(x) = \frac{1}{h_1(h_1 + h_2)/2} (f_0 - f_1) + \frac{1}{h_2(h_1 + h_2)/2} (f_2 - f_1) =$$

$$= \left(\frac{f_2 - f_1}{h_2} - \frac{f_1 - f_0}{h_1} \right) \frac{1}{(h_1 + h_2)/2} = \frac{2}{h_1 + h_2} [f_x(x_1) - f_{\bar{x}}(x_1)]. \quad (14.20)$$

Упражнение 14.2. Доказать, что соотношение (14.20) при $h_1 \neq h_2$ аппроксимирует производную $f''(x_1)$ с погрешностью не выше $O(h_1 + h_2)$. При какой гладкости $f(x)$ такая погрешность достигается?

Для построения формул численного дифференцирования можно использовать не только интерполяционные многочлены Лагранжа, но и интерполяционные многочлены Эрмита. Такие формулы полезны, когда в узлах заданы не только значения

функции, но и значения производных, а производные нужно знать в других точках. Формула численного дифференцирования и в этом случае выглядит аналогично (14.19). Если

$$f(x) = H_m(x) + R_m(x),$$

где $H_m(x)$ — интерполяционный многочлен Эрмита, а $R_m(x)$ — остаточный член, то

$$f^{(k)}(x) \approx H_m^{(k)}(x), \quad 1 \leq k \leq m. \quad (14.21)$$

Пример 14.6. Пусть интерполяционный многочлен Эрмита имеет степень три и написан по значениям функции и ее первой производной в двух узлах $x_0 = 0$ и $x_1 = h$, т.е.

$$H_3(x) = p_{00}(x)f_0 + p_{01}(x)f'_0 + p_{10}(x)f_1 + p_{11}(x)f'_1,$$

где

$$p_{00}(x) = \frac{(2x+h)(x-h)^2}{h^3}, \quad p_{01}(x) = \frac{x(x-h)^2}{h^2},$$

$$p_{10}(x) = \frac{x^2(3h-2x)}{h^3}, \quad p_{11}(x) = \frac{x^2(x-h)}{h^2}.$$

Тогда формула для приближенного нахождения первой производной в точке x примет вид

$$H'_3(x) = p'_{00}(x)f_0 + p'_{01}(x)f'_0 + p'_{10}(x)f_1 + p'_{11}(x)f'_1.$$

Если принять во внимание, что

$$p'_{00}(x) = 6x(x-h)/h^3, \quad p'_{01}(x) = \frac{3x^2 - 4hx + h^2}{h^2},$$

$$p'_{10}(x) = -6x(x-h)/h^3, \quad p'_{11}(x) = \frac{3x^2 - 2hx}{h^2},$$

то, например, $H'_3(0) = f'_0$, а

$$H'_3(h/2) = \frac{3}{2} \frac{f_1 - f_0}{h} - \frac{f'_0 + f'_1}{4}.$$

Далее, для приближенного вычисления второй производной в точке x имеем формулу

$$H''_3(x) = p''_{00}(x)f_0 + p''_{01}(x)f'_0 + p''_{10}(x)f_1 + p''_{11}(x)f'_1.$$

Принимая во внимание, что

$$p''_{00}(x) = 6(2x-h)/h^3, \quad p''_{01}(x) = 2(3x-2h)/h^2,$$

$$p''_{10}(x) = 6(-2x+h)/h^3, \quad p''_{11}(x) = 2(3x-h)/h^2,$$

Находим, например,

$$H''_3(0) = \frac{6(f_1 - f_0)}{h^2} - \frac{4}{h}f'_0 - \frac{2}{h}f'_1.$$

Элементарные вычисления показывают, что

$$6(f_1 - f_0)/h^2 - 4f'_0/h - 2f'_1/h = f''_0 - \frac{h^2}{12}f^{IV}(\xi).$$

14.4 О корректности численного дифференцирования

В формулах численного дифференцирования линейные комбинации значений функции $f(x)$ в узлах x_i делятся на h^m , где m — порядок вычисляемой производной. Поскольку сами значения функции, как правило, задаются или вычисляются не точно, то при малых h неустранимые погрешности оказывают существенное влияние на точность численного дифференцирования.

Пусть δ_i — величина погрешности, с которой вычисляется значение функции $f(x)$ в узле x_i , т.е. вычисляемое приближенное значение есть

$$\tilde{f}_i = f_i + \delta_i.$$

Будем предполагать, что $|\delta_i| \leq \delta$.

Пусть для приближенного вычисления первой производной используется формула (14.8). Тогда, с учетом (14.7),

$$\begin{aligned} f'(x_i) &\approx \frac{\tilde{f}_{i+1} - \tilde{f}_{i-1}}{2h} = \frac{f_{i+1} + \delta_{i+1} - f_{i-1} - \delta_{i-1}}{2h} = \\ &= f'(x_i) + \frac{h^2}{6} f'''(\xi_i) + (\delta_{i+1} - \delta_{i-1})/(2h). \end{aligned}$$

Отсюда находим, что для полной погрешности этой формулы $\varepsilon_1 = (\tilde{f}_{i+1} - \tilde{f}_{i-1})/(2h) - f'(x_i)$ справедлива оценка

$$|\varepsilon_1| \leq \frac{h^2}{6} M_3 + \delta/h, \quad (14.22)$$

где $M_3 = \max_{x \in [x_{i+1}, x_{i-1}]} |f'''|$. Из этой оценки следует, что при уменьшении h полная погрешность убывает только до определенного предела, после чего начинает расти. Если, например, δ сравнима с h , то *мы не можем найти приближенное значение производной*, ибо погрешность будет $O(1)$. Чтобы вычисленное значение можно было рассматривать как приближенное значение производной, нужно, чтобы h было много больше δ . Наивысшую точность мы получим при том h , при котором правая часть (14.22) достигает минимума по h . Указанное значение

$$h = h_1 = \sqrt[3]{3\delta/M_3}.$$

При этом

$$\varepsilon_1 = \frac{3}{2} \left(\frac{M_3}{3} \right)^{1/3} \delta^{2/3}.$$

Если при тех же предположениях о точности вычисления значений f_i воспользоваться формулой из левой части (14.10), дающее приближенное значение второй производной, то полная погрешность $\varepsilon_2 = (\tilde{f}_{i+1} - 2\tilde{f}_i + \tilde{f}_{i-1})/(h^2) - f''(x_i)$ оценится так

$$|\varepsilon_2| \leq \frac{h^2}{12} M_4 + \frac{4\delta}{h^2},$$

где $M_4 = \max_{x_{i-1} \leq x \leq x_{i+1}} |f^{IV}|$. При этом оптимальное значение $h = h_2 = 2(3\delta/M_4)^{1/4}$, а $\varepsilon_2 = 2\sqrt{M_4/3} \delta^{1/2}$.

Следует, однако, заметить, что предельная точность при приближенном вычислении производных не всегда ниже, чем точность, с которой задана сама функция. Пусть, например, $\tilde{f}_i = f_i + \delta v_i$, где v_i — некоторая "гладкая функция", т.е. такая, что, например, $|v_{x,i}| \leq M$. Тогда для формулы (14.8) полная погрешность будет оцениваться следующим образом:

$$|\varepsilon_1| \leq \frac{h^2}{6} M_3 + M\delta$$

и, если $h/\sqrt{\delta} = O(1)$, то $|\varepsilon_1| = O(\delta)$.

В.Б. Андреев

ЧИСЛЕННЫЕ МЕТОДЫ

Часть II

Глава IV

Методы решения задачи Коши для обыкновенных дифференциальных уравнений

§ 15

Постановка задачи и первые примеры

15.1 Введение. Задача Коши

Рассмотрим задачу Коши для обыкновенного дифференциального уравнения первого порядка

$$\frac{du}{dt} = f(t, u), \quad t > 0, \quad u(0) = u_0. \quad (15.1)$$

Из курса дифференциальных уравнений известно, что для однозначной разрешимости задачи (15.1) в некоторой окрестности точки $t = 0$ достаточно, чтобы функция $f(t, u)$ была непрерывна в окрестности точки $(0, u_0)$ и удовлетворяла условию Липшица по второму аргументу. Известны примеры, иллюстрирующие отсутствие решения задачи (15.1) или его неединственность при нарушении указанных условий. Мы всегда будем предполагать, что решение задачи (15.1) существует и единственно. Для дальнейшего нам даже придется предполагать, что искомое решение достаточно гладкое.

15.2 Примеры численных методов

Приведем несколько простейших численных методов решения задачи (15.1). Для этого введем на полуоси $t \geq 0$ равномерную сетку, т.е. множество точек (которые назовем узлами)

$$\omega = \{t_n = n\tau, \quad n = 0, 1, \dots; \quad \tau > 0\}$$

и будем искать приближенное решение задачи (15.1) в узлах ω . Величину τ будем называть шагом сетки ω . Договоримся приближенное решение в узле t_n обозначать той же буквой, что и решение задачи (15.1), но с индексом n внизу: u_n . Тем самым, мы отказываемся от часто используемого обозначения $u(t_n) = u_n$; теперь $u(t_n)$ — значение точного решения в узле t_n , а u_n — значение приближенного решения в этом узле, и, вообще говоря, $u(t_n) \neq u_n$. Наоборот, $u_n - u(t_n)$ представляет собой погрешность численного метода в узле t_n , которую нам предстоит оценивать. Данное соглашение не представляется наилучшим, однако остановимся на нем.

Для построения численных методов проинтегрируем уравнение (15.1) от t_n до t_{n+1}

$$u(t_{n+1}) - u(t_n) = \int_{t_n}^{t_{n+1}} f(t, u(t)) dt \quad (15.2)$$

и заменим приближенно интеграл в правой части этой формулы какой-либо квадратурной формулой. Здесь мы рассмотрим четыре таких формулы.

Построенная в курсе "Введение в численные методы" квадратурная формула прямоугольников представляет интеграл произведением длины отрезка интегрирования и значения подынтегральной функции в середине этого отрезка

$$\int_a^b \varphi(x) dx \approx |b - a| \varphi\left(\frac{a + b}{2}\right). \quad (15.3)$$

Эта квадратурная формула точна на многочленах первой степени, и при малых $|b - a|$ ее погрешность есть $O(|b - a|^3)$.

Наряду с этой квадратурной формулой, которую мы впредь будем называть формулой *центральных прямоугольников*, можно ввести так называемые формулы *левых* и *правых прямоугольников*. Первая из них состоит в представлении интеграла произведением длины отрезка интегрирования и значения подынтегральной функции в левом конце отрезка

$$\int_a^b \varphi(x) dx \approx |b - a| \varphi(a), \quad (15.4)$$

а вторая — произведением длины отрезка интегрирования и значения подынтегральной функции в правом конце отрезка

$$\int_a^b \varphi(x) dx \approx |b - a| \varphi(b). \quad (15.5)$$

Обе эти формулы точны только на многочленах нулевой степени и имеют погрешность $O(|b - a|^2)$.

а) **Метод Эйлера.** Заменим интеграл в (15.2) формулой левых прямоугольников (15.4). В результате получим приближенное равенство

$$u(t_{n+1}) - u(t_n) \approx \tau f(t_n, u(t_n)). \quad (15.6)$$

Определим приближенное решение задачи (15.1) как такую сеточную функцию, заданную на ω , которое превращает соотношение (15.6) в равенство. Разделив полученное равенство на τ , будем иметь

$$\frac{u_{n+1} - u_n}{\tau} = f(t_n, u_n), \quad n = 0, 1, \dots, \quad u_0 = u(0). \quad (15.7)$$

Соотношение (15.7) позволяет рекуррентным образом найти приближенное решение во всех узлах. Численный метод решения задачи (15.1), реализуемый формулами (15.7), называется *методом Эйлера*.

б) **Неявный метод Эйлера.** Заменяем теперь интеграл в (15.2) формулой правых прямоугольников (15.5). Для отыскания приближенного решения получим уравнения

$$\frac{u_{n+1} - u_n}{\tau} = f(t_{n+1}, u_{n+1}), \quad n = 0, 1, \dots, \quad u_0 = u(0). \quad (15.8)$$

Соотношения (15.8) коренным образом отличаются от соотношений (15.7): для отыскания приближенного решения u_{n+1} теперь нужно решать нелинейные уравнения

$$u_{n+1} - \tau f(t_{n+1}, u_{n+1}) = u_n.$$

Метод (15.8) называется *неявным методом Эйлера*. С точки зрения простоты вычислений он сильно уступает обычному методу Эйлера (15.7). Как будет показано позже, по точности оба метода сравнимы. Еще позже будет установлена существенно бóльшая устойчивость метода (15.8) по сравнению с (15.7).

в) **Метод Рунге.** Заменяем интеграл в (15.2) формулой центральных прямоугольников (15.3)

$$u(t_{n+1}) - u(t_n) \approx \tau f(t_{n+1/2}, u(t_{n+1/2})). \quad (15.9)$$

Использованный нами ранее прием получения численного метода путем превращения приближенного равенства в точное за счет замены $u(t_n)$ на u_n здесь напрямую не проходит: в приближенном равенстве фигурирует значение f при u в точке $t_{n+1/2}$, которая не является узловой. Если же мы все же воспользуемся этим приемом и введем *промежуточное* значение приближенного решения в точке $t_{n+1/2}$, то нам потребуется дополнительное уравнение для определения приближенного решения в точке $t_{n+1/2}$. Обозначим промежуточное значение приближенного решения через $u_{n+1/2}$. Тогда из (15.9)

$$\frac{u_{n+1} - u_n}{\tau} = f(t_{n+1/2}, u_{n+1/2}), \quad (15.10)$$

а для отыскания $u_{n+1/2}$ напомним, например, соотношение Эйлера (15.7)

$$\frac{u_{n+1/2} - u_n}{\tau/2} = f(t_n, u_n). \quad (15.11)$$

Итак, в методе (15.10), (15.11) вычисление нового приближенного значения искомого решения u_{n+1} осуществляется поэтапно. Сначала находится промежуточное значение $u_{n+1/2}$ по формуле (15.11), а затем и само u_{n+1} из (15.10). Вычисления по обеим формулам явные. *Метод* (15.10), (15.11) был предложен немецким математиком *Рунге* и носит его имя. Будет показано, что точность метода (15.10), (15.11) выше, чем точность методов (15.7) и (15.8); для вычисления интеграла все же использована более точная квадратурная формула.

Замечание 15.1. В некоторых учебниках по численным методам метод (15.10), (15.11) называется методом предиктор-корректор (предсказывающе корректирующим).

г) **Метод трапеций.** Наконец, заменим интеграл в (15.2) формулой трапеций. В результате получим

$$\frac{u_{n+1} - u_n}{\tau} = \frac{f(t_n, u_n) + f(t_{n+1}, u_{n+1})}{2}, \quad n = 0, 1, \dots, \quad u_0 = u(0). \quad (15.12)$$

Как и в случае неявного метода Эйлера (15.8), реализация метода (15.12) требует решения нелинейного уравнения

$$u_{n+1} - \frac{\tau}{2} f(t_{n+1}, u_{n+1}) = F(u_n).$$

Будет показано, что точность метода (15.12) сравнима с точностью метода Рунге (15.10), (15.11), а по устойчивости он значительно превосходит последний и в этом отношении близок к неявному методу Эйлера (15.8). Метод (15.12) иногда называют *методом трапеций*.

15.3 Аппроксимация.

Определение 15.1. Сеточная функция

$$z_n = u_n - u(t_n), \quad n = 1, 2, \dots$$

называется погрешностью решения.

Замечание 15.2. Погрешность решения определена только в узлах основной сетки ω , но не в промежуточных узлах.

Выведем уравнение, которому удовлетворяет погрешность решения в методе Эйлера (15.7). Подставив $u_n = z_n + u(t_n)$ в (15.7), получим

$$\frac{z_{n+1} - z_n}{\tau} + \frac{u(t_{n+1}) - u(t_n)}{\tau} = f(t_n, u(t_n) + z_n). \quad (15.13)$$

Преобразуем правую часть этого соотношения путем разложения по формуле Тейлора

$$f(t_n, u(t_n) + z_n) = f(t_n, u(t_n)) + z_n \frac{\partial f}{\partial u}(t_n, \tilde{u}),$$

где

$$\tilde{u} = u(t_n) + \theta z_n, \quad 0 < \theta < 1.$$

Подставляя это разложение в (15.13) и преобразовывая, найдем, что

$$\frac{z_{n+1} - z_n}{\tau} = \frac{\partial f}{\partial u}(t_n, \tilde{u}) z_n + \psi_n, \quad (15.14)$$

где

$$\psi_n = f(t_n, u(t_n)) - \frac{u(t_{n+1}) - u(t_n)}{\tau}. \quad (15.15)$$

Искомое уравнение получено.

Определение 15.2. Сеточная функция ψ_n , задаваемая соотношением (15.15), называется *погрешностью аппроксимации* дифференциального уравнения (15.1) уравнением (15.7).

Замечание 15.3. Погрешность аппроксимации представляет собой разность между правой и левой частями уравнения, определяющего численный метод, если туда вместо приближенного решения подставить точное.

Оценим погрешность аппроксимации метода Эйлера. Используя формулу Тейлора и принимая во внимание уравнение (15.1), в предположении непрерывности второй производной $u(t)$, из (15.15) будем иметь

$$\begin{aligned}\psi_n &= f(t_n, u(t_n)) - \frac{u(t_n) + \tau u'(t_n) + \frac{\tau^2}{2} u''(t_n + \theta\tau) - u(t_n)}{\tau} = \\ &= [f(t_n, u(t_n)) - u'(t_n)] + \frac{\tau}{2} u''(t_n + \theta\tau) = O(\tau).\end{aligned}$$

Тем самым, метод Эйлера имеет первый порядок аппроксимации.

Упражнение 15.1. Исследовать погрешности аппроксимации методов (15.8) и (15.12).

Указание. Для упрощения выкладок разложение по формуле Тейлора в методе (15.8) вести в точке t_{n+1} , а в методе (15.12) — в точке $t_{n+1/2}$.

§ 16

Методы Рунге-Кутты

16.1 Общая концепция

Численные методы решения уравнения

$$\frac{du}{dt} = f(t, u), \quad t > 0, \quad u(0) = u_0 \quad (16.1)$$

и систем таких уравнений, наиболее широко используемые в вычислительной практике, делятся на два больших класса: многошаговые методы и методы типа Рунге-Кутты. Все приведенные в качестве примеров численные методы относятся к методам Рунге-Кутты, хотя некоторые из них могут трактоваться и как многошаговые (одношаговые).

Сейчас мы опишем общую концепцию методов Рунге-Кутты. Для этого вновь обратимся к интегральному соотношению (15.2), на основе которого мы строили изложенные выше методы. Но прежде сделаем одно допущение относительно уравнения (16.1), которое в дальнейшем существенно облегчит нам жизнь. Будем предполагать, что правая часть f этого уравнения не зависит явным образом от t , т.е. $f \equiv f(u)$ и, следовательно,

$$\frac{du}{dt} = f(u), \quad t > 0, \quad u(0) = u_0. \quad (16.1')$$

Сделанное допущение не является ограничением, ибо все численные методы, построенные для одного уравнения, допускают очевидное распространение на случай системы, т.е., вообще говоря, u можно считать вектором. Если же f зависит явным образом от t , то, обозначив например, $t = u_0(t)$ и объявив $u_0(t)$ новой неизвестной, удовлетворяющей уравнению

$$u_0'(t) = 1, \quad u_0(0) = 0,$$

мы сведем задачу к ранее оговоренному случаю.

Итак, пусть $f = f(u)$. Перепишем для этого случая интегральное соотношение (15.2)

$$u(t_{n+1}) - u(t_n) = \int_{t_n}^{t_{n+1}} f(u) dt. \quad (16.2)$$

Сделаем в интеграле (16.2) замену переменной интегрирования, полагая

$$(t - t_n)/\tau = \theta. \quad (16.3)$$

Эта замена переводит отрезок $[t_n, t_{n+1}]$ в $[0, 1]$ так, что

$$u(t_{n+1}) - u(t_n) = \tau \int_0^1 f(\hat{u}(\theta)) d\theta, \quad (16.4)$$

где

$$\hat{u}(\theta) = u(t(\theta)).$$

Пусть

$$0 \leq \theta_1 < \theta_2 < \dots < \theta_s \leq 1 \quad (16.5)$$

суть узлы, а b_1, b_2, \dots, b_s — веса некоторой квадратурной формулы, аппроксимирующей интеграл $\int_0^1 \varphi(\theta) d\theta$. Используя эту формулу для аппроксимации интеграла в (16.4), будем иметь

$$u(t_{n+1}) - u(t_n) \approx \tau \sum_{i=1}^s b_i f(\hat{u}(\theta_i)). \quad (16.6)$$

Чтобы получить из этого соотношения численный метод, нужно точные значения искомого решения заменить на приближенные, а приближенное равенство — наоборот на точное. Но прежде мы должны ввести дополнительные обозначения. Будем обозначать значение приближенного решения в точке t , отвечающей узлу квадратурной формулы θ_i ($t = t_n + \tau\theta_i$) через Y_i . Тогда искомое уравнение примет вид

$$u_{n+1} = u_n + \tau \sum_{i=1}^s b_i f(Y_i). \quad (16.7)$$

Чтобы получить уравнение для определения Y_i , проинтегрируем (16.1') от t_n до $t_n + \tau\theta_i$ и сделаем замену (16.3)

$$\hat{u}(\theta_i) - u(t_n) = \int_{t_n}^{t_n + \tau\theta_i} f(u(t)) dt = \tau \int_0^{\theta_i} f(\hat{u}(\theta)) d\theta.$$

Заменим и здесь интеграл квадратурной формулой с теми же узлами (16.5). Эта квадратурная формула будет несколько своеобразной, ибо не все ее узлы будут лежать

на отрезке интегрирования. Разумеется, ее веса, вообще говоря, должны быть отличны от b_j и даже быть своими для каждого i . Пусть

$$Y_i = u_n + \tau \sum_{j=1}^s a_{ij} f(Y_j), \quad i = \overline{1, s}. \tag{16.8}$$

Соотношения (16.7), (16.8) полностью определяют численный метод.

Итак, для того, чтобы найти приближенное решение u_{n+1} (когда u_n уже найдено), сначала нужно решить, вообще говоря, нелинейную систему (16.8) и определить Y_i , $i = \overline{1, s}$, которые затем следует подставить в (16.7).

Определение 16.1. Метод (16.7), (16.8) называется *s-этапным методом Рунге-Кутты*.

Этот метод принято записывать таблицей его коэффициентов, которая называется таблицей Бутчера

$$\begin{array}{c|cccc}
 c_1 & a_{11} & a_{12} & \dots & a_{1s} \\
 c_2 & a_{21} & a_{22} & \dots & a_{2s} \\
 & \dots & \dots & \dots & \dots \\
 c_s & a_{s1} & a_{s2} & \dots & a_{ss} \\
 \hline
 & b_1 & b_2 & \dots & b_s
 \end{array}
 \quad c_i = \sum_{j=1}^s a_{ij}. \tag{16.9}$$

Замечание 16.1. Поскольку b_i суть весовые коэффициенты квадратурной формулы для интеграла по единичному отрезку, то $\sum_{i=1}^s b_i = 1$. Из аналогичных соображений

$$c_i = \sum_{j=1}^s a_{ij} = \theta_i.$$

Определение 16.2. Если в таблице Бутчера (16.9) коэффициенты $a_{ij} = 0$ при $j \geq i$, то метод (16.7), (16.8) называется явным *s-этапным методом Рунге-Кутты*.

Определение 16.3. Если $a_{ij} = 0$ при $i > j$ и хотя бы один $a_{ii} \neq 0$, то метод (16.7), (16.8) называется диагонально неявным.

Во всех остальных случаях мы говорим о неявных методах Рунге-Кутты.

Коэффициенты в таблице Бутчера (16.9) при заданных ограничениях выбираются из соображений максимальной точности численного метода.

16.2 Одноэтапные методы Рунге-Кутты

Исследуем одноэтапные ($s = 1$) методы Рунге-Кутты. При $s = 1$ соотношения (16.8), (16.7) принимают вид

$$Y_1 = u_n + \tau a_{11} f(Y_1), \tag{16.10}$$

$$u_{n+1} = u_n + \tau b_1 f(Y_1) \tag{16.11}$$

Из соображений аппроксимации (квадратурная формула должна быть точной по крайней мере на const) находим, что $b_1 = 1$. Если теперь положить $a_{11} = 0$, то метод будет явным, причем $Y_1 = u_n$, а (16.11) можно переписать в виде

$$\frac{u_{n+1} - u_n}{\tau} = f(u_n).$$

Мы получили метод Эйлера. Тем самым, метод Эйлера есть *явный одноэтапный метод Рунге-Кутты*.

Если взять $a_{11} = 1$, то метод (16.10), (16.11) будет неявным. При этом правые части (16.10) и (16.11) совпадают и приводят к соотношению $Y_1 = u_{n+1}$. В этом случае система (16.10), (16.11) преобразуется к виду

$$\frac{u_{n+1} - u_n}{\tau} = f(u_{n+1}).$$

Это неявный метод Эйлера (15.8). Он также является одноэтапным методом Рунге-Кутты.

Исследуем теперь наиболее целесообразный выбор параметров b_1 и a_{11} с точки зрения минимизации погрешности аппроксимации. Чтобы найти погрешность аппроксимации, перепишем уравнение (16.11) в виде

$$\frac{u_{n+1} - u_n}{\tau} = b_1 f(Y_1) \quad (16.12)$$

(ср. с (15.7), (15.8), (15.10) и (15.12)), а решение уравнения (16.10) обозначим через $Y_1(u_n)$. Если, как и выше, $z_n = u_n - u(t_n)$, то

$$\frac{z_{n+1} - z_n}{\tau} = b_1 f(Y_1(u(t_n) + z_n)) - \frac{u(t_{n+1}) - u(t_n)}{\tau}.$$

И снова, раскладывая первое слагаемое правой части по формуле Тейлора, находим, что

$$\begin{aligned} \frac{z_{n+1} - z_n}{\tau} &= b_1 \left[f(Y_1(u(t_n))) + \frac{\partial f}{\partial u}(\tilde{u})z_n \right] - \frac{u(t_{n+1}) - u(t_n)}{\tau} = \\ &= b_1 \frac{\partial f}{\partial Y_1} \frac{\partial Y_1}{\partial u}(\tilde{u})z_n + \psi_n, \end{aligned}$$

где

$$\psi_n = b_1 f(Y_1(u(t_n))) - \frac{u(t_{n+1}) - u(t_n)}{\tau} \quad (16.13)$$

— погрешность аппроксимации, а $Y_1(u(t_n))$ — решение уравнения (16.10) с $u(t_n)$ вместо u_n , т.е.

$$Y_1(u(t_n)) = u(t_n) + \tau a_{11} f(Y_1(u(t_n))). \quad (16.14)$$

Замечание 16.2. Погрешность аппроксимации (16.13) представляет собой разность между правой и левой частями уравнения (16.12), если туда вместо приближенного решения подставить точное (ср. с замечанием 15.3).

Разложим погрешность аппроксимации (16.13) по степеням τ . Имеем

$$\psi_n = b_1 \left[f(Y_1) \Big|_{\tau=0} + \tau \frac{df(Y_1)}{d\tau} \Big|_{\tau=0} + \frac{\tau^2}{2} \frac{d^2 f}{d\tau^2} \right] - \left[u'(t_n) + \frac{\tau}{2} u''(t_n) + \frac{\tau^2}{6} \tilde{u}''' \right].$$

Из (16.14) находим, что $Y_1|_{\tau=0} = u(t_n)$ и, следовательно,

$$f(Y_1) \Big|_{\tau=0} = f(u(t_n)).$$

Снова с использованием (16.14)

$$\frac{df(Y_1)}{d\tau} \Big|_{\tau=0} = \frac{df}{dY_1} \frac{dY_1}{d\tau} \Big|_{\tau=0} = \frac{df}{du}(u(t_n)) a_{11} f(u(t_n)),$$

а из уравнения (16.1')

$$u'(t_n) = f(u(t_n)), \quad u''(t_n) = \frac{df}{dt}(u(t_n)) = \frac{df}{du}(u(t_n)) \frac{du}{dt}(t_n) = \frac{df}{du} f.$$

Поэтому

$$\psi_n = (b_1 - 1)f(u(t_n)) + \tau \left[b_1 a_{11} - \frac{1}{2} \right] f(u(t_n)) \frac{df}{du}(u(t_n)) + O(\tau^2).$$

Тем самым, для того, чтобы погрешность аппроксимации была $O(\tau^2)$, необходимо и достаточно, чтобы выполнялись условия

$$b_1 = 1, \quad a_{11} b_1 = 1/2. \quad (16.15)$$

Отсюда находим

$$b_1 = 1, \quad a_{11} = 1/2$$

и, следовательно, неявный одноэтапный метод Рунге-Кутты

$$\begin{aligned} Y_1 &= u_n + \frac{\tau}{2} f(Y_1), \\ u_{n+1} &= u_n + \tau f(Y_1) \end{aligned} \quad (16.16)$$

имеет второй порядок аппроксимации.

Замечание 16.3. Из первого уравнения (16.16) следует, что момент времени, на который Y_1 приближает $u(t)$, есть $t + \tau/2$, ибо для задачи $u' = 1$, $u(0) = 0$, имеющей решение $u = t$, $Y_1 = u_n + \tau/2 = t_n + \tau/2$.

Соотношения (16.16) можно преобразовать. Исключив $f(Y_1)$, найдем, что $u_{n+1} = 2Y_1 - u_n$. Выражая отсюда Y_1 и подставляя его во второе уравнение (16.16), получим

$$u_{n+1} = u_n + \tau f\left(\frac{u_{n+1} + u_n}{2}\right).$$

Замечание 16.4. Метод (16.16) очень сильно напоминает метод Рунге (15.10), (15.11). Отличие между ними состоит в том, что здесь промежуточное значение находится по неявной формуле, а в методе Рунге по явной формуле (15.11). Метод (16.16), как мы уже сказали, является одноэтапным (неявным) методом Рунге-Кутты, а метод (15.10), (15.11) — двухэтапным (явным) методом. Подчеркнем, что слову *этап* здесь мы придаем четкий математический смысл.

16.3 Методы третьего порядка аппроксимации

Выясним ограничения на коэффициенты (16.9), обеспечивающие третий порядок аппроксимации s -этапного метода Рунге-Кутты. Для этого нужно исследовать погрешность аппроксимации

$$\psi_n := \psi_n(\tau) := \sum_{i=1}^s b_i f(Y_i(u(t_n))) - \frac{u(t_{n+1}) - u(t_n)}{\tau},$$

где

$$Y_i(u(t_n)) = u(t_n) + \tau \sum_{j=1}^s a_{ij} f(Y_j(u(t_n))) =: Y_i(u(t_n); \tau). \quad (16.17)$$

Раскладывая $\psi_n(\tau)$ по τ до третьего порядка, будем иметь

$$\begin{aligned} \psi_n(\tau) = & \sum_{i=1}^s b_i \left[f(Y_i) \Big|_{\tau=0} + \tau \frac{df(Y_i)}{d\tau} \Big|_{\tau=0} + \frac{\tau^2}{2} \frac{d^2 f(Y_i)}{d\tau^2} \Big|_{\tau=0} + O(\tau^3) \right] - \\ & - \left[u'(t_n) + \frac{\tau}{2} u''(t_n) + \frac{\tau^2}{6} u'''(t_n) + O(\tau^3) \right]. \end{aligned} \quad (16.18)$$

Поскольку $f(Y_i(u(t_n)))$ есть сложная функция τ , то вычислим сначала производные по τ функции $Y_i(u(t_n); \tau)$ при $\tau = 0$. Из (16.17) с учетом (16.9), находим, что

$$\begin{aligned} Y_i \Big|_{\tau=0} &= u(t_n), \\ Y_i' \Big|_{\tau=0} &= \frac{dY_i}{d\tau} \Big|_{\tau=0} = \left[\sum_{j=1}^s a_{ij} f(Y_j) + \tau \sum_{j=1}^s a_{ij} \frac{df}{dY_j} Y_j' \right] \Big|_{\tau=0} = f(u(t_n)) c_i, \\ Y_i'' \Big|_{\tau=0} &= \left[2 \sum_{j=1}^s a_{ij} \frac{df}{dY_j} Y_j' + \tau \sum_{j=1}^s a_{ij} \frac{d^2 f}{dY_j^2} (Y_j')^2 + \tau \sum_{j=1}^s a_{ij} \frac{df}{dY_j} Y_j'' \right] \Big|_{\tau=0} = 2f(u(t_n)) \frac{df}{du} \sum_{j=1}^s a_{ij} c_j. \end{aligned}$$

Теперь можно найти производные f :

$$\begin{aligned} f(Y_i(u(t_n))) \Big|_{\tau=0} &= f(u(t_n)), \\ \frac{df(Y_i)}{d\tau} \Big|_{\tau=0} &= \frac{df}{dY_i} Y_i' \Big|_{\tau=0} = f(u(t_n)) \frac{df}{du} c_i, \\ \frac{d^2 f(Y_i)}{d\tau^2} \Big|_{\tau=0} &= \left[\frac{d^2 f}{dY_i^2} (Y_i')^2 + \frac{df}{dY_i} Y_i'' \right] \Big|_{\tau=0} = f^2(u(t_n)) \frac{d^2 f}{du^2} c_i^2 + 2f(u(t_n)) \left(\frac{df}{du} \right)^2 \sum_{j=1}^s a_{ij} c_j. \end{aligned}$$

Далее, из (16.1')

$$u' = f, \quad u'' = \frac{df}{du} u' = f' f, \quad u''' = f'' u' f + (f')^2 u' = f'' f^2 + (f')^2 f$$

и, следовательно,

$$\frac{u(t_{n+1}) - u(t_n)}{\tau} = f + \frac{\tau}{2} f' f + \frac{\tau^2}{6} (f'' f^2 + (f')^2 f) + O(\tau^3).$$

Подставляя теперь найденные разложения в (16.18), будем иметь

$$\begin{aligned} \psi_n &= \sum_{i=1}^s b_i \left[f + \tau f f' c_i + \frac{\tau^2}{2} \left(f^2 f'' c_i^2 + 2f f'^2 \sum_{j=1}^s a_{ij} c_j \right) \right] - \\ &- \left[f + \frac{\tau}{2} f f' + \frac{\tau^2}{6} (f^2 f'' + f f'^2) \right] + O(\tau^3). \end{aligned} \quad (16.19)$$

Отсюда, приравнявая коэффициенты при одинаковых степенях τ , находим, что условия третьего порядка аппроксимации суть

$$\begin{aligned} \sum_{i=1}^s b_i &= 1, \\ \sum_{i=1}^s b_i c_i &= \frac{1}{2}, \end{aligned} \quad (16.20)$$

$$\begin{aligned} \sum_{i=1}^s b_i c_i^2 &= \frac{1}{3}, \\ \sum_{i,j=1}^s b_i a_{ij} c_j &= \frac{1}{6}. \end{aligned} \quad (16.21)$$

При этом (16.20) суть условия второго порядка аппроксимации.

Замечание 16.5. Чтобы иметь условия четвертого порядка аппроксимации, к условиям (16.20), (16.21) нужно добавить следующие условия:

$$\begin{aligned}
 \sum_{i=1}^s b_i c_i^3 &= \frac{1}{4}, \\
 \sum_{i,j=1}^s b_i c_i a_{ij} c_j &= \frac{1}{8}, \\
 \sum_{i,j=1}^s b_i a_{ij} c_j^2 &= \frac{1}{12}, \\
 \sum_{i,j,k=1}^s b_i a_{ij} a_{jk} c_k &= \frac{1}{24}.
 \end{aligned} \tag{16.22}$$

Замечание 16.6. Условия (16.20) с учетом замечания 16.1 можно трактовать как условия точности квадратурной формулы из (16.6) на линейных функциях.¹ Добавление к этим условиям первого из соотношений (16.21), а затем и первого из соотношений (16.22) на указанную квадратурную формулу накладывает дополнительные условия точности на квадратичных и кубических функциях.

Упражнение 16.1. Показать, что метод трапеций (15.12) является неявным двухэтапным методом Рунге-Кутты второго порядка аппроксимации. (Найти все b_i , a_{ij} и показать невыполнение хотя бы одно из условий (16.21))

Ответ:

$$\begin{array}{c|cc}
 0 & 0 & 0 \\
 1 & 1/2 & 1/2 \\
 \hline
 & 1/2 & 1/2
 \end{array} \quad \left(\frac{1}{2}0^2 + \frac{1}{2}1^2 \right) \neq \frac{1}{3}.$$

Упражнение 16.2. Показать, что метод Рунге (15.10), (15.11) является явным двухэтапным методом Рунге-Кутты второго порядка.

Ответ:

$$\begin{array}{c|cc}
 1/2 & 1/2 & \\
 \hline
 & 0 & 1
 \end{array} \quad \left[0 \cdot 0 + 1 \cdot \frac{1}{4} \right] \neq \frac{1}{3}.$$

16.4 Двухэтапные неявные методы третьего порядка

Положим в (16.20), (16.21) параметр $s = 2$. В результате система примет вид

$$\begin{aligned}
 b_1 + b_2 &= 1, \\
 c_1 b_1 + c_2 b_2 &= 1/2, \\
 c_1^2 b_1 + c_2^2 b_2 &= 1/3, \\
 b_1(a_{11}c_1 + a_{12}c_2) + b_2(a_{21}c_1 + a_{22}c_2) &= 1/6.
 \end{aligned} \tag{16.23}$$

¹Ведь $c_i = \theta_i$, т.е. координата переменной интегрирования в i -ом узле.

Эта система содержит четыре уравнения и шесть неизвестных (Если не считать c_1 и c_2 , задаваемые (16.9)). Поэтому, вообще говоря, два из этих неизвестных должны остаться свободными, а остальные выразиться через них. Система (16.23) нелинейная, и нет регулярных способов ее решения. Укажем один путь, приводящий к решению этой системы.

Для отыскания решения системы (16.23) предположим сначала, что неизвестные c_1 и c_2 найдены, и рассмотрим первые три уравнения (16.23) как систему линейных уравнений относительно b_1 и b_2 . Поскольку эта система переопределена, то для ее разрешимости необходимо обращение в нуль определителя расширенной матрицы

$$\begin{aligned} \begin{vmatrix} 1 & 1 & 1 \\ c_1 & c_2 & 1/2 \\ c_1^2 & c_2^2 & 1/3 \end{vmatrix} &= \frac{1}{3}c_2 + \frac{1}{2}c_1^2 + c_1c_2^2 - c_1^2c_2 - \frac{1}{2}c_2^2 - \frac{1}{3}c_1 = \\ &= \frac{1}{3}(c_2 - c_1) - \frac{1}{2}(c_2 + c_1)(c_2 - c_1) + c_1c_2(c_2 - c_1) = \\ &= (c_2 - c_1) \left[\frac{1}{3} - \frac{c_1 + c_2}{2} + c_1c_2 \right] = 0. \end{aligned} \quad (16.24)$$

Проанализируем это соотношение. Если бы $c_1 = c_2$, то последнее уравнение (16.23) приняло бы вид

$$c_1^2b_1 + c_2^2b_2 = 1/6,$$

что противоречит третьему уравнению (16.23), и поэтому

$$c_1 - c_2 \neq 0. \quad (16.25)$$

Тем самым, из (16.24) следует, что

$$2 - 3(c_1 + c_2) + 6c_1c_2 = 0$$

или

$$(3 - 6c_1)c_2 = 2 - 3c_1.$$

Поскольку $c_1 = 1/2$ не удовлетворяет этому уравнению, то

$$c_1 \neq 1/2 \quad (16.26)$$

и можно найти

$$c_2 = \frac{2 - 3c_1}{3 - 6c_1}. \quad (16.27)$$

Разрешим теперь первые два уравнения (16.23) относительно b_1 и b_2 при помощи правила Крамера. Будем иметь

$$\Delta = \begin{vmatrix} 1 & 1 \\ c_1 & c_2 \end{vmatrix} = c_2 - c_1 \neq 0, \quad \Delta_1 = \begin{vmatrix} 1 & 1 \\ 1/2 & c_2 \end{vmatrix} = c_2 - 1/2, \quad \Delta_2 = \begin{vmatrix} 1 & 1 \\ c_1 & 1/2 \end{vmatrix} = 1/2 - c_1$$

и, следовательно,

$$b_1 = \frac{c_2 - 1/2}{c_2 - c_1} = \frac{1}{4(3c_1^2 - 3c_1 + 1)}, \quad b_2 = \frac{1/2 - c_1}{c_2 - c_1}. \quad (16.28)$$

Из (16.27), (16.28) следует, что c_1 можно принять за параметр. В качестве второго параметра возьмем a_{12} . Тогда

$$a_{11} = c_1 - a_{12}. \quad (16.29)$$

Поскольку

$$a_{21} = c_2 - a_{22}, \quad (16.30)$$

то, подставляя эти выражения для a_{11} и a_{21} в последнее из уравнений (16.23), получим

$$b_1[(c_1 - a_{12})c_1 + a_{12}c_2] + b_2[(c_2 - a_{22})c_1 + a_{22}c_2] = 1/6.$$

Принимая во внимание (16.28), второе из уравнений (16.23) и разрешая полученное соотношение относительно a_{22} , будем иметь

$$a_{22} = \frac{1/6 - c_1/2 - a_{12}(c_2 - 1/2)}{1/2 - c_1} = \frac{(1 - 3c_1)(1 - 2c_1) - a_{12}}{3(1 - 2c_1)^2}. \quad (16.31)$$

Соотношения (16.27), (16.28), (16.29), (16.30), (16.31) задают двухпараметрическое семейство неявных двухэтапных методов Рунге-Кутты третьего порядка.

Если положить, например,

$$a_{12} = 0, \quad c_1 \equiv a_{11} = a_{22}, \quad (16.32)$$

то для c_1 из (16.30) получим квадратное уравнение $6c_1^2 - 6c_1 + 1$ с корнями $c_1 = \gamma = \frac{3 \pm \sqrt{3}}{6}$. Таблица Бутчера для этого метода имеет вид

$$\begin{array}{c|cc} \theta_1 = \gamma & \gamma & 0 \\ \theta_2 = 1 - \gamma & 1 - 2\gamma & \gamma \\ \hline & 1/2 & 1/2 \end{array} \quad \gamma = \frac{3 \pm \sqrt{3}}{6}. \quad (16.33)$$

Упражнение 16.3. Доказать, что метод (16.33) есть метод (16.27)-(16.32).

16.5 Явные двухэтапные методы

В силу определения для явного двухэтапного метода $a_{11} = a_{12} = a_{22} = 0$ и лишь $a_{21} \neq 0$. (При $a_{21} = 0$ мы получим явный одноэтапный метод.) Поскольку двухэтапные методы третьего порядка имеют лишь два свободных параметра, а мы задали три, то рассчитывать на третий порядок у явных двухэтапных методов, вообще говоря, не приходится. Мы покажем, что так оно и есть.

Принимая a_{21} за параметр, из условий второго порядка аппроксимации (16.20), которые в нашем случае принимают вид

$$b_1 + b_2 = 1, \quad a_{21}b_2 = 1/2,$$

находим

$$b_1 = \left(1 - \frac{1}{2a_{21}}\right), \quad b_2 = \frac{1}{2a_{21}}.$$

Тем самым, явные двухэтапные методы Рунге-Кутты второго проядка образуют однопараметрическое семейство.

Далее, поскольку в рассматриваемом случае наряду с a_{11} , a_{12} , a_{22} и $c_1 = 0$, то левая часть четвертого из условий (16.23) обращается в нуль и следовательно это условие выполненным быть не может. Мы доказали, что явных двухэтапных методов третьего порядка не существует.

Если положить, например, $a_{21} = 1$, то получим метод Хойна

$$\begin{array}{l} Y_1 = u_n, \quad Y_2 = u_n + \tau f(Y_1), \\ u_{n+1} = u_n + \frac{\tau}{2}[f(Y_1) + f(Y_2)]. \end{array} \quad \begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & 1/2 & 1/2 \end{array}$$

Упражнение 16.4. Выписать все построенные методы второго порядка.

16.6 Двухэтапный метод четвертого порядка

Коэффициенты метода четвертого порядка должны удовлетворять еще четырем условиям (16.22). Хотя у двухэтапного метода третьего порядка осталось только два параметра, существует единственный двухэтапный метод четвертого порядка. Его коэффициенты суть

$$\begin{array}{c|cc} 1/2 - \sqrt{3}/6 & 1/4 & 1/4 - \sqrt{3}/6 \\ 1/2 + \sqrt{3}/6 & 1/4 + \sqrt{3}/6 & 1/4 \\ \hline & 1/2 & 1/2 \end{array} \quad (16.34)$$

Замечание 16.7. Если бы f не зависела от u ($u' = f(t)$, $u_{n+1} = u_n + \tau \int_0^1 \hat{f}(\theta) d\theta$), то двухэтапный метод четвертого порядка получился бы только если квадратура в (16.7) была квадратурой Гаусса, т.е. $b_1 = b_2 = 1/2$, а θ_1 и θ_2 — сдвинутые на $[0, 1]$ нули полинома Лежандра второй степени $P_2(x) = \frac{1}{2}(3x^2 - 1)$. Корнями этого полинома являются числа $x_{1,2} = \pm 1/\sqrt{3}$. Делая линейную замену, переводящую отрезок $[-1, 1]$ в отрезок $[0, 1]$, находим, что узлы квадратуры Гаусса на $[0, 1]$ суть $\theta_{1,2} = 1/2 \mp \sqrt{3}/6$, как в (16.34). Остальные коэффициенты получаются, если проинтегрировать по $[0, \theta_1]$ и $[0, \theta_2]$ весовые функции интерполяционного полинома Лагранжа с гауссовыми узлами.

Упражнение 16.5. Доказать, что в (16.34)

$$a_{ij} = \int_0^{\theta_j} p_i(\theta) d\theta,$$

где $p_i(\theta)$ — линейная функция (интерполюант по двум узлам) такая, что $p_i(\theta_i) = 1$, $p_i(\theta_j) = 0$ при $i \neq j$. Убедиться в выполнении условий (16.20)-(16.22).

16.7 Явные трехэтапные методы Рунге-Кутты третьего порядка

Рассмотрим более подробно явные трехэтапные методы. В силу определения

$$a_{11} = a_{12} = a_{13} = a_{22} = a_{23} = a_{33} = 0,$$

и указанные методы задаются таблицей

$$\begin{array}{c|ccc} c_2 & a_{21} & & \\ c_3 & a_{31} & a_{32} & \\ \hline & b_1 & b_2 & b_3 \end{array}$$

Условия третьего порядка аппроксимации (16.20), (16.21) в рассматриваемом случае принимают вид

$$\begin{aligned} b_1 + b_2 + b_3 &= 1, \\ b_2 c_2 + b_3 c_3 &= 1/2, \\ b_2 c_2^2 + b_3 c_3^2 &= 1/3, \\ b_3 a_{32} c_2 &= 1/6. \end{aligned} \tag{16.35}$$

Эта система имеет два однопараметрических семейства решений и одно двухпараметрическое. Найдем их. Будем рассматривать второе и третье уравнения системы (16.35) как линейную систему относительно b_2 и b_3 . Эта система может быть как вырожденной (и это приводит к двум однопараметрическим семействам решений), так и невырожденной (двухпараметрическое семейство).

Пусть эта система вырождена, т.е.

$$\begin{vmatrix} c_2 & c_3 \\ c_2^2 & c_3^2 \end{vmatrix} = c_2 c_3 (c_3 - c_2) = 0. \tag{16.36}$$

В силу последнего из уравнений (16.35) $c_2 \neq 0$, и поэтому либо

$$c_3 = 0 \tag{16.37}$$

либо

$$c_2 = c_3. \tag{16.38}$$

i) Пусть сначала имеет место (16.37). Тогда второе и третье уравнения (16.35) принимают вид

$$c_2 b_2 = 1/2, \quad c_2^2 b_2 = 1/3$$

и, следовательно,

$$c_2 = 2/3, \quad b_2 = 3/4.$$

Если теперь $b_3 = b$ принять за параметр, то из последнего уравнения (16.35) находим

$$a_{32} = \frac{1}{4b}.$$

Поскольку $c_3 = 0$, то

$$a_{31} = -a_{32} = -\frac{1}{4b}$$

и, наконец, из первого уравнения (16.35)

$$b_1 = \frac{1}{4} - b.$$

Собирая найденные значения, получим таблицу

$$\begin{array}{c|cc} 2/3 & 2/3 & \\ 0 & -(4b)^{-1} & (4b)^{-1} \\ \hline & 1/4 - b & 3/4 \quad b \end{array} \quad (16.39)$$

ii) Теперь пусть имеет место (16.38). Снова из второго и третьего уравнений (16.35) находим, что

$$b_2 + b_3 = \frac{1}{2c_2} = \frac{1}{3c_2^2},$$

т.е.

$$c_3 = c_2 = 2/3, \quad b_2 = 3/4 - b,$$

где $b = b_3$ — параметр. Из последнего уравнения (16.35)

$$a_{32} = \frac{1}{4b},$$

а из первого уравнения

$$b_1 = 1/4.$$

Таблица рассматриваемого метода имеет вид

$$\begin{array}{c|ccc} 2/3 & 2/3 & & \\ 2/3 & \frac{2}{3} - \frac{1}{4b} & \frac{1}{4b} & \\ \hline & 1/4 & 3/4 - b & b \end{array} \quad (16.40)$$

iii) Если соотношение (16.36) места не имеет, то из второго и третьего уравнений (16.35) находим

$$b_2 = \frac{c_3/2 - 1/3}{c_2(c_3 - c_2)}, \quad b_3 = \frac{1/3 - c_2/2}{c_3(c_3 - c_2)}, \quad (16.41)$$

Привлекая первое уравнение (16.35), найдем, что

$$b_1 = 1 - \frac{3(c_2 + c_3) - 2}{6c_2c_3}, \quad (16.42)$$

а из четвертого

$$a_{32} = \frac{c_3(c_3 - c_2)}{c_2(2 - 3c_2)}. \quad (16.43)$$

Еще раз напомним, что в методе (16.41)-(16.43) параметры c_2 и c_3 удовлетворяют условию

$$c_2c_3(c_3 - c_2)(c_2 - 2/3) \neq 0.$$

Среди явных трехэтапных методов Рунге-Кутты третьего порядка в силу исторических причин наибольшей популярностью пользуется следующий метод из семейства (16.41)-(16.43)

$$\begin{array}{c|ccc} 1/2 & 1/2 & & \\ 1 & -1 & 2 & \\ \hline & 1/6 & 2/3 & 1/6 \end{array}. \quad (16.44)$$

16.8 Более общие методы Рунге-Кутты

Теорема 16.1. *Не существует явного s -этапного метода Рунге-Кутты порядка p , если $p > s$.*

Для $s = 1, 2$ эта теорема нами фактически доказана.

Теорема 16.2. *При $s \geq 5$ не существует явного s -этапного метода Рунге-Кутты порядка $p = s$ (1963г.). При $s \geq 8$ не существует явного s -этапного метода Рунге-Кутты порядка $p = s - 1$ (1965г.). При $s \geq 10$ — $p = s - 2$ (1985г.).*

Замечание 16.8. При $s = 6$ существует явный метод Рунге-Кутты порядка 5. При $s = 7$ существует явный метод порядка 6. Наивысший порядок, фактически достигнутый для явно построенных явных методов Рунге-Кутты равен 10. При этом число этапов равно 17. (Этот результат занесен в книгу рекордов Гиннеса).

Теорема 16.3. *При любом s существует единственный неявный метод Рунге-Кутты порядка $p = 2s$.*

Замечание 16.9. Для оптимального метода порядка $p = 2s$ узлы θ_j и веса b_j суть узлы и веса квадратурной формулы Гаусса, а

$$a_{ij} = \int_0^{\theta_j} p_i(\theta) d\theta,$$

где $p_j(\theta)$ — многочлен степени s такой, что $p_i(\theta_i) = 1$, $p_i(\theta_j) = 0$ при $i \neq j$.

В вычислительной практике широко используется следующий явный 4-этапный метод Рунге-Кутты 4-го порядка

$$\begin{array}{c|cccc} 1/2 & 1/2 & & & \\ 1/2 & 0 & 1/2 & & \\ 1 & 0 & 0 & 1 & \\ \hline & 1/6 & 1/3 & 1/3 & 1/6. \end{array}$$

Замечание 16.10. До недавнего времени в методах Рунге-Кутты (явных) вместо переменных Y_j фигурировали $k_j = f(Y_j)$. Поэтому, вместо (16.7), (16.8) писали

$$\begin{aligned} k_i &= f\left(u_n + \tau \sum_{j=1}^{s-1} a_{ij} k_j\right), \\ u_{n+1} &= u_n + \tau \sum_{j=1}^s b_j k_j. \end{aligned} \tag{16.45}$$

16.9 Сходимость методов Рунге-Кутты

Установим оценку погрешности приближенного решения, получаемого при помощи того или иного метода Рунге-Кутты.

Если, как и раньше,

$$z_n = u_n - u(t_n),$$

то из (16.7) находим, что z_{n+1} удовлетворяет уравнению

$$\frac{z_{n+1} - z_n}{\tau} = \sum_{i=1}^s b_i \frac{df(Y_i(u))}{du} \Bigg|_{u=u(t_n)+\sigma_i z_n} z_n + \psi_n. \tag{16.46}$$

Прежде чем оценивать z_{n+1} , оценим коэффициент при z_n в правой части (16.46). Будем при этом предполагать, что

$$\max_{|u| < \infty} \left| \frac{df(u)}{du} \right| \leq L. \tag{16.47}$$

Тогда

$$\max_{|u|<\infty} \left| \frac{df(Y_j(u))}{du} \right| = \max_{|u|<\infty} \left| \frac{df}{dY_j} \frac{dY_j}{du} \right| \leq L \max_{|u|<\infty} \left| \frac{dY_j}{du} \right|. \quad (16.48)$$

Оценим $|dY_j/du|$. Из (16.8) с u вместо u_n

$$\frac{dY_i}{du} = 1 + \tau \sum_{j=1}^s a_{ij} \frac{df}{dY_j} \frac{dY_j}{du}.$$

Пусть

$$\max_{|u|<\infty} \left| \frac{dY_{i_0}}{du} \right| = \max_{1 \leq j \leq s} \max_{|u|<\infty} \left| \frac{dY_j}{du} \right|.$$

Тогда с учетом (16.47)

$$\max_{|u|<\infty} \left| \frac{dY_{i_0}}{du} \right| \leq 1 + \tau \sum_{j=1}^s |a_{i_0 j}| L \max_{|u|<\infty} \left| \frac{dY_j}{du} \right| \leq 1 + \tau a s L \max_{|u|<\infty} \left| \frac{dY_{i_0}}{du} \right|,$$

где

$$a = \max_{ij} |a_{ij}|, \quad (16.49)$$

и, следовательно,

$$(1 - \tau a s L) \max_{|u|<\infty} \left| \frac{dY_{i_0}}{du} \right| \leq 1.$$

Будем предполагать, что

$$\tau \leq \frac{1}{2asL}, \quad \text{т.е.} \quad 1 - \tau a s L \geq \frac{1}{2}. \quad (16.50)$$

Тогда

$$\left| \frac{dY_j}{du} \right| \leq 2, \quad j = 1, \dots, s. \quad (16.51)$$

Для простоты будем предполагать, что коэффициенты b_j неотрицательны. Поскольку их сумма равна единице, то с учетом (16.48), (16.51)

$$\left| \sum_{j=1}^s b_j \frac{df(Y_j(u))}{du} \right| \leq 2L.$$

Принимая во внимание эту оценку, из (16.46) находим, что

$$|z_{n+1}| \leq (1 + 2\tau L)|z_n| + \tau|\psi_n|.$$

Разрешим эти неравенства. Поскольку $z_0 = u_0 - u(0) = 0$, то

$$\begin{aligned} |z_1| &\leq \tau|\psi_0|, \\ |z_2| &\leq (1 + 2\tau L)|z_1| + \tau|\psi_1|, \\ |z_3| &\leq (1 + 2\tau L)|z_2| + \tau|\psi_2|, \\ &\dots\dots\dots \\ |z_n| &\leq (1 + 2\tau L)|z_{n-1}| + \tau|\psi_{n-1}| \end{aligned}$$

Подставим теперь оценку $|z_1|$ в правую часть оценки $|z_2|$, а полученную оценку $|z_2|$ в свою очередь в правую часть оценки $|z_3|$ и т.д. Получим

$$\begin{aligned}
 |z_n| &\leq \sum_{j=0}^{n-1} \tau(1 + 2\tau L)^{n-1-j} |\psi_j| \leq (1 + 2\tau L)^n \sum_{j=0}^{n-1} \tau |\psi_j| \leq e^{2\tau L T / \tau} T \max_j |\psi_j| = \\
 &= e^{2LT} T \max_j |\psi_j|.
 \end{aligned} \tag{16.52}$$

Из (16.52) следует

Теорема 16.4. *Если метод Рунге-Кутты (16.7), (16.8) аппроксимирует исходное уравнение (16.1') с порядком p , то при $\tau \rightarrow 0$ он сходится с тем же порядком.*

§ 17

Линейные многошаговые методы

При изучении методов Рунге-Кутты, используемых при решении задачи Коши

$$\frac{du}{dt} = f(u), \quad t > 0, \quad u(0) = u_0, \quad (17.1)$$

мы не обращали особого внимания на задание начальных условий, ибо это совершенно тривиальная процедура: для того, чтобы начал работать любой из рассмотренных нами методов Рунге-Кутты, нужно задать $u_0 = u(0)$, т.е. так же как и для дифференциального уравнения. Обусловлено это тем, что в каждом уравнении связаны между собой значения приближенного решения в двух соседних узлах сетки (не считая промежуточных значений). Другой класс методов составляют так называемые многошаговые методы, в которых уравнения связывают значения приближенного решения в нескольких соседних узлах.

17.1 Методы Адамса

Наиболее известными из многошаговых методов и наиболее старыми являются методы Адамса. Опишем эти методы на примере уравнения (17.1). Вновь будем предполагать, что на отрезке интегрирования введена равномерная сетка с шагом τ , а уравнение (17.1) проинтегрировано по отрезку между узлами t_n и t_{n+1}

$$u(t_{n+1}) - u(t_n) = \int_{t_n}^{t_{n+1}} f(u(t)) dt. \quad (17.2)$$

Заменим подынтегральную функцию в (17.2) интерполяционным многочленом Лагранжа по некоторым узлам сетки ω (а не по промежуточным узлам, как это было в методах Рунге-Кутты (!)). В зависимости от того, участвует ли узел t_{n+1} в интерполяции $f(u(t))$ или нет, различают неявные и явные методы Адамса.

а) **Явные методы Адамса.** Предположим, что $u(t)$ известна в k узлах сетки ω

$$t_n, t_{n-1}, \dots, t_{n+1-k}. \quad (17.3)$$

Построим по этим узлам интерполяционный многочлен Лагранжа степени $k - 1$ для подынтегральной функции $f(u(t))$ из (17.2)

$$f(u(t)) \approx L_{k-1}(t) = \sum_{j=1}^k p_j(t) f(u(t_{n+1-j})), \quad (17.4)$$

где, как обычно,

$$p_j(t) = \prod_{\substack{i=1 \\ i \neq j}}^k \frac{t - t_{n+1-i}}{t_{n+1-j} - t_{n+1-i}} \quad (17.5)$$

суть весовые функции интерполяционного полинома (многочлены степени $(k - 1)$), обращающиеся в нуль при $t = t_{n+1-i}$, $i = \overline{1, j-1, j+1, k}$ и в единицу при $t = t_{n+1-j}$. Подставляя (17.4), (17.5) в (17.2) и заменяя приближенное равенство на точное, получим следующее уравнение для определения приближенного решения

$$u_{n+1} - u_n = \tau \sum_{j=1}^k b_j f(u_{n+1-j}), \quad (17.6)$$

где

$$b_j = \frac{1}{\tau} \int_{t_n}^{t_{n+1}} p_j(t) dt = \int_0^1 \hat{p}_j(\theta) d\theta = \int_0^1 \prod_{\substack{i=1 \\ i \neq j}}^k \frac{\theta - 1 + i}{i - j} d\theta. \quad (17.7)$$

Определение 17.1. Численный метод (17.6), (17.7) называется явным k -шаговым методом Адамса (иногда его называют методом Адамса-Бэшфорта).

Примеры. 1°. $k = 1$.

$$p_1(t) = \hat{p}_1(\theta) = 1, \quad b_1 = 1.$$

2°. $k = 2$.

$$\begin{aligned} \hat{p}_1(\theta) &= \theta + 1, & b_1 &= 3/2, \\ \hat{p}_2(\theta) &= -\theta, & b_2 &= -1/2. \end{aligned}$$

3°. $k = 3$.

$$\begin{aligned} \hat{p}_1(\theta) &= \frac{1}{2}(\theta + 1)(\theta + 2), & b_1 &= \frac{23}{12}, \\ \hat{p}_2(\theta) &= -\theta(\theta + 2), & b_2 &= -\frac{4}{3}, \\ \hat{p}_3(\theta) &= \frac{\theta(\theta + 1)}{2}, & b_3 &= \frac{5}{12}. \end{aligned}$$

Выпишем уравнения (17.6) для этих частных случаев

$$\begin{aligned} u_{n+1} &= u_n + \tau f(u_n), \\ u_{n+1} &= u_n + \tau \left[\frac{3}{2}f(u_n) - \frac{1}{2}f(u_{n-1}) \right], \\ u_{n+1} &= u_n + \tau \left[\frac{23}{12}f(u_n) - \frac{16}{12}f(u_{n-1}) + \frac{5}{12}f(u_{n-2}) \right]. \end{aligned} \quad (17.8)$$

Упражнение 17.1. Построить явный 4-х-шаговый метод Адамса (17.6).

Ответ.

$$u_{n+1} = u_n + \tau \left[\frac{55}{24}f(u_n) - \frac{59}{24}f(u_{n-1}) + \frac{37}{24}f(u_{n-2}) - \frac{9}{24}f(u_{n-3}) \right].$$

Замечание 17.1. Очевидно, что первое из уравнений (17.8) определяет исследованный нами ранее метод Эйлера. Тем самым, метод Эйлера может быть отнесен как к методам Рунге-Кутты, так и к методам Адамса.

Формулы (17.6) получены при интегрировании в пределах от t_n до t_{n+1} , в то время как узлы интерполяции располагались на отрезке $[t_{n+1-k}, t_n]$, т.е. вне интервала интегрирования (Для подынтегральной функции использовалась экстраполяция). В связи с этим явные методы Адамса иногда называют экстраполяционными методами.

б) **Неявные методы Адамса.** Можно построить и неявные методы Адамса. Для этого к узлам интерполяции (17.3) нужно добавить еще узел t_{n+1} . В этом случае интерполяционный многочлен (степени k) примет вид

$$L_k(t) = \sum_{j=0}^k p_j(t)f(u(t_{n+1-j})), \quad (17.9)$$

а соответствующим ему уравнением будет уравнение

$$u_{n+1} - u_n = \tau \sum_{j=0}^k b_j f(u_{n+1-j}), \quad (17.10)$$

где (ср. с (17.7))

$$b_j = \frac{1}{\tau} \int_{t_n}^{t_{n+1}} p_j(t) dt = \int_0^1 \prod_{\substack{i=0 \\ i \neq j}}^k \frac{\theta - 1 + i}{i - j} d\theta. \quad (17.11)$$

Определение 17.2. Численный метод (17.10), (17.11) называется неявным k -шаговым методом Адамса (Иногда его называют методом Адамса-Мултона).

Примеры. 4°. $k = 0$.

$$\hat{p}_0(\theta) = 1, \quad b_0 = 1.$$

5°. $k = 1$.

$$\begin{aligned} \hat{p}_0(\theta) &= \theta, & b_0 &= 1/2, \\ \hat{p}_1(\theta) &= -\theta + 1, & b_1 &= 1/2. \end{aligned}$$

6°. $k = 2$.

$$\begin{aligned} \hat{p}_0(\theta) &= \frac{1}{2}\theta(\theta + 1), & b_0 &= \frac{5}{12}, \\ \hat{p}_1(\theta) &= -(\theta^2 - 1), & b_1 &= \frac{2}{3}, \\ \hat{p}_2(\theta) &= \frac{1}{2}\theta(\theta - 1), & b_2 &= -\frac{1}{12}. \end{aligned}$$

Напишем уравнения (17.10) для этих частных случаев

$$\begin{aligned} u_{n+1} &= u_n + \tau f(u_{n+1}), \\ u_{n+1} &= u_n + \frac{\tau}{2} [f(u_{n+1}) + f(u_n)], \\ u_{n+1} &= u_n + \tau \left[\frac{5}{12}f(u_{n+1}) + \frac{8}{12}f(u_n) - \frac{1}{12}f(u_{n-1}) \right]. \end{aligned} \tag{17.12}$$

Упражнение 17.2. Построить неявный 3-х-шаговый метод Адамса (17.10), (17.11).

Ответ.

$$u_{n+1} = u_n + \tau \left[\frac{9}{24}f(u_{n+1}) + \frac{19}{24}f(u_n) - \frac{5}{24}f(u_{n-1}) + \frac{1}{24}f(u_{n-2}) \right].$$

Замечание 17.2. Очевидно, что первое из уравнений (17.12), отвечающее $k = 0$, является неявным методом Эйлера, а второе уравнение, отвечающее $k = 1$, — методом трапеций. Тем самым, эти одношаговые неявные методы Адамса являются и методами Рунге-Кутты.

17.2 Формулы дифференцирования назад

Во всех предыдущих случаях, как при построении методов Рунге-Кутты, так и при построении методов Адамса, мы получали численные методы путем интегрирования уравнения (17.1) и замены подынтегральной функции $f(u)$ в (17.2) интерполяционным многочленом или замены интеграла квадратурной формулой. А можно поступать и иначе: интерполяционным многочленом заменить $u(t)$. Тогда для построения численного метода нужно будет выражение интерполяционного многочлена подставить в

(17.1). Чтобы получился численный метод, точка t_{n+1} должна быть в числе узлов интерполяции. Пусть

$$u(t) \approx L_k(t) = \sum_{j=0}^k p_j(t)u(t_{n+1-j}), \quad (17.13)$$

где

$$p_j(t) = \prod_{\substack{i=0 \\ i \neq j}}^k \frac{t - t_{n+1-i}}{t_{n+1-j} - t_{n+1-i}}.$$

Подставляя (17.13) в (17.1), получим приближенное равенство

$$\sum_{j=0}^k p'_j(t)u(t_{n+1-j}) \approx f\left(\sum_{j=0}^k p_j(t)u(t_{n+1-j})\right).$$

Превратим его в точное равенство в каком-либо узле. В результате получим уравнение для определения приближенного решения. Рассмотрим случай, когда указанным узлом является t_{n+1} . Будем иметь

$$\sum_{j=0}^k p'_j(t_{n+1})u_{n+1-j} = f(u_{n+1}).$$

Как и раньше, сделаем локальную замену переменной $(t - t_n)/\tau = \theta$. Тогда

$$p'_j(t) = \frac{dp_j(t)}{dt} = \frac{1}{\tau} \frac{d\hat{p}_j(\theta)}{d\theta} = \frac{1}{\tau} \hat{p}'_j(\theta),$$

где

$$p_j(t) = \hat{p}_j(\theta) = \prod_{\substack{i=0 \\ i \neq j}}^k \frac{\theta - 1 + i}{i - j},$$

и полученный метод принимает вид

$$\sum_{j=0}^k \hat{p}'_j(1)u_{n+1-j} = \tau f(u_{n+1}). \quad (17.14)$$

Определение 17.3. Численные методы (17.14) называются формулами дифференцирования назад.

Примеры. 7°. $k = 1$.

$$\begin{aligned} \hat{p}_0(\theta) &= \theta, & \hat{p}'_0(1) &= 1, \\ \hat{p}_1(\theta) &= -\theta + 1, & \hat{p}'_1(1) &= -1. \end{aligned}$$

8°. $k = 2$.

$$\begin{aligned}\hat{p}_0(\theta) &= \frac{1}{2}\theta(\theta + 1), & \hat{p}'_0(1) &= \frac{3}{2}, \\ \hat{p}_1(\theta) &= 1 - \theta^2, & \hat{p}'_1(1) &= -2, \\ \hat{p}_2(\theta) &= \frac{1}{2}\theta(\theta - 1), & \hat{p}'_2(1) &= \frac{1}{2}.\end{aligned}$$

Выпишем уравнения (17.14) для этих случаев

$$u_{n+1} - u_n = \tau f(u_{n+1}), \quad (17.15)$$

$$\left(\frac{3}{2}u_{n+1} - 2u_n + \frac{1}{2}u_{n-1}\right) = \tau f(u_{n+1}). \quad (17.16)$$

Упражнение 17.3. Построить формулу (17.14), отвечающую $k = 3$.

Ответ.

$$\left(\frac{11}{6}u_{n+1} - 3u_n + \frac{3}{2}u_{n-1} - \frac{1}{3}u_{n-2}\right) = \tau f(u_{n+1}).$$

17.3 Общие линейные многошаговые методы

Методы Адамса, явные и неявные, и формулы дифференцирования назад являются частными случаями формулы

$$\sum_{j=0}^k \alpha_j u_{n-j} = \tau \sum_{j=0}^k \beta_j f(u_{n-j}), \quad (17.17)$$

где α_j и β_j — действительные числа. (Обратим внимание на то, что в этой формуле вместо нового неизвестного u_{n+1} фигурирует u_n). Будет предполагать, что

$$\alpha_0 \neq 0, \quad |\alpha_k| + |\beta_k| \neq 0. \quad (17.18)$$

Первое из условий (17.18) обеспечивает разрешимость неявного ($\beta_0 \neq 0$) уравнения (17.17) по крайней мере, для достаточно малого шага τ . Второе из условий (17.18) всегда можно считать выполненным, уменьшив при необходимости k .

Определение 17.4. Формула (17.17) называется линейным многошаговым (k -шаговым) методом.

Метод является явным, если $\beta_0 = 0$, и неявным в противном случае.

Чтобы линейный многошаговый метод (17.17) можно было использовать для численного решения задачи (17.1), необходимо, чтобы уравнение (17.17) аппроксимировало уравнение (17.1).

Определение 17.5. Величина

$$\psi_n = \sum_{j=0}^k \beta_j f(u(t_{n-j})) - \frac{1}{\tau} \sum_{j=0}^k \alpha_j u(t_{n-j}) \quad (17.19)$$

называется погрешностью аппроксимации метода (17.17).

Выясним вопрос о порядке погрешности аппроксимации метода (17.17) при $\tau \rightarrow 0$.

Теорема 17.1. *Многошаговый метод (17.17) имеет погрешность аппроксимации порядка $p \leq 2R$ тогда и только тогда, когда выполняются следующие условия*

$$\sum_{j=0}^k \alpha_j = 0, \quad \sum_{j=0}^k (\alpha_j j^q + q \beta_j j^{q-1}) = 0, \quad q = 1, \dots, p. \quad (17.20)$$

Доказательство. Разложим $u(t)$ по формуле Тейлора в точке t_n :

$$u(t) = \sum_{q=0}^p \frac{(t-t_n)^q}{q!} u^{(q)}(t_n) + O((t-t_n)^{p+1}). \quad (17.21)$$

Так как $f(u) = u'(t)$, то, дифференцируя (17.21), получим

$$f(u(t)) = \sum_{q=0}^p q \frac{(t-t_n)^{q-1}}{q!} u^{(q)}(t_n) + O((t-t_n)^p). \quad (17.22)$$

Подставляя теперь разложения (17.21), (17.22) при $t = t_{n-j}$ в (17.19), будем иметь

$$\begin{aligned} \psi_n &= \sum_{j=0}^k \beta_j \sum_{q=0}^p q \frac{(-j\tau)^{q-1}}{q!} u^{(q)}(t_n) - \\ &- \frac{1}{\tau} \sum_{j=0}^k \alpha_j \sum_{q=0}^p \frac{(-j\tau)^q}{q!} u^{(q)}(t_n) + O(\tau^p) = \\ &= \sum_{q=0}^p \frac{(-\tau)^{q-1}}{q!} u^{(q)}(t_n) \sum_{j=0}^k [\beta_j q j^{q-1} + \alpha_j j^q] + O(\tau^p). \end{aligned}$$

Приравнивая нулю коэффициенты при τ^{q-1} для $q = 0, 1, \dots, p$, получим (17.20). Теорема доказана.

Замечание 17.3. Решение уравнения (17.17) не изменится, если его умножить на какое-либо число, отличное от нуля. Это означает, что его коэффициенты определяются с точностью до множителя (до мультипликативной постоянной). Чтобы устранить этот произвол, пронормируем их, полагая, например,

$$\sum_{j=0}^k \beta_j = 1. \quad (17.23)$$

Замечание 17.4. Из (17.20), (17.23) имеем $(p+2)$ уравнения для $2(k+1)$ коэффициентов метода (17.17). Тем самым, максимальный порядок аппроксимации линейного k -шагового метода есть $p = 2k$.

17.4 Погрешность аппроксимации методов Адамса

Исследуем вопрос о порядке погрешности аппроксимации методов Адамса. Для этого перепишем сначала явный метод Адамса (17.6), (17.7) в виде (17.17), т.е. заменим $n+1$ на n :

$$u_n - u_{n-1} = \tau \sum_{j=1}^k b_j f(u_{n-j}).$$

Сравнивая это соотношение с (17.17), находим, что

$$\alpha_0 = 1, \alpha_1 = -1, \alpha_2 = \dots = \alpha_k = 0, \beta_0 = 0, b_j = \beta_j, j = \overline{1, k}.$$

Определим, для каких дифференциальных уравнений явные методы Адамса теоретически дают точное решение в узлах сетки. Это произойдет в том случае, когда интерполяционный многочлен $L_{k-1}(t)$, определяющий явный метод Адамса, совпадает с $f(u)$ или с $f(t, u)$. Пусть $f(t, u(t)) = f(t)$, т.е. f не зависит от u и является многочленом степени не выше $k-1$. Тогда $f(t)$ совпадает со своим интерполяционным многочленом $L_{k-1}(t)$, и явный метод Адамса точен для уравнений

$$u' = qt^{q-1}, \quad q = 0, \dots, k.$$

Это означает, что погрешность аппроксимации (17.19) на решениях этих уравнений равна нулю. Подставляя решения этих уравнений $u = t^q$ в (17.19) при $n = 0$, получим

$$\psi_0 = \sum_{j=0}^k \left[\beta_j q (-\tau j)^{q-1} - \frac{1}{\tau} \alpha_j (-\tau j)^q \right] = 0, \quad q = 0, \dots, k,$$

что совпадает с первыми $(k+1)$ уравнениями (17.20). Тем самым, мы доказали, что явный k -шаговый метод Адамса имеет порядок погрешности аппроксимации не ниже k . Можно показать, что его порядок аппроксимации в точности равен k .

Упражнение 17.4. Доказать, что порядок аппроксимации неявного k -шагового метода Адамса не ниже $k+1$.

Упражнение 17.5. Доказать, что порядок аппроксимации k -шаговой формулы дифференцирования назад не ниже k .

17.5 Поучительный пример

Построим двухшаговый явный метод максимального порядка аппроксимации. Согласно ранее сказанному, порядок аппроксимации этого метода должен быть равен трем. Из (17.20), (17.23) имеем

$$\begin{aligned}\alpha_0 + \alpha_1 + \alpha_2 &= 0, \\ \alpha_1 + 2\alpha_2 &= -(\beta_0 + \beta_1 + \beta_2), \\ \alpha_1 + 4\alpha_2 &= -2(\beta_1 + 2\beta_2), \\ \alpha_1 + 8\alpha_2 &= -3(\beta_1 + 4\beta_2), \\ \beta_0 + \beta_1 + \beta_2 &= 1, \\ \beta_0 &= 0.\end{aligned}$$

Разрешая эту линейную систему, находим, что

$$\alpha_0 = \frac{1}{6}, \quad \alpha_1 = \frac{2}{3}, \quad \alpha_2 = -\frac{5}{6}, \quad \beta_1 = \frac{2}{3}, \quad \beta_2 = \frac{1}{3}.$$

Тем самым, метод (17.17) приобретает вид

$$\left(\frac{1}{6}u_n + \frac{4}{6}u_{n-1} - \frac{5}{6}u_{n-2}\right) = \tau \left[\frac{2}{3}f_{n-1} + \frac{1}{3}f_{n-2}\right]. \quad (17.24)$$

Применим этот метод к решению уравнения (17.1) с $f(u) = \lambda u$, где $\lambda = \text{const}$. Будем при этом предполагать, что начальное значение $u_0 = 1$. В этом случае задача (17.1) примет вид

$$u'(t) = \lambda u, \quad u(0) = 1, \quad (17.25)$$

а ее решением будет функция

$$u(t) = e^{\lambda t}. \quad (17.26)$$

Отвечающий (17.25) метод (17.17) можно записать так

$$\sum_{j=0}^k (\alpha_j - \tau \lambda \beta_j) u_{n-j} = 0, \quad (17.27)$$

а применительно к методу (17.24)

$$\frac{1}{6}u_n + \left(\frac{4}{6} - \frac{2}{3}\tau\lambda\right)u_{n-1} + \left(-\frac{5}{6} - \frac{1}{3}\tau\lambda\right)u_{n-2} = 0. \quad (17.28)$$

Это есть линейное однородное разностное уравнение второго порядка с постоянными коэффициентами (см. §6). Найдем его решение. Для этого нужно написать характеристическое уравнение, отвечающее разностному уравнению (17.28), и найти его корни. Искомое характеристическое уравнение имеет вид

$$q^2 + 4(1 - \tau\lambda)q - (5 + 2\tau\lambda) = 0, \quad (17.29)$$

а его корни суть

$$\begin{aligned} q_1 &= -2 + 2\tau\lambda + \sqrt{9 - 6\tau\lambda + 4\tau^2\lambda^2} = 1 + \tau\lambda + O(\tau^2\lambda^2), \\ q_2 &= -2 + 2\tau\lambda - \sqrt{9 - 6\tau\lambda + 4\tau^2\lambda^2} = -5 + O(\tau\lambda). \end{aligned} \quad (17.30)$$

Упражнение 17.6. Доказать, что $q_1 - e^{\tau\lambda} = O(\tau^4\lambda^4)$.

Поскольку корни (17.30) характеристического уравнения различны, то общее решение разностного уравнения (17.28) имеет вид

$$u_n = c_1 q_1^n + c_2 q_2^n, \quad (17.31)$$

где c_1 и c_2 — произвольные постоянные.

Рассматриваемый нами метод (17.24) является двухшаговым, и одного начального условия

$$u_0 = 1 \quad (17.32)$$

для его реализации недостаточно. Поскольку точное решение нам известно, то не будем ломать голову над тем, как задать недостающее начальное условие при $n = 1$, а просто положим

$$u_1 = u(t_1) = e^{\tau\lambda}. \quad (17.33)$$

Потребуем, чтобы решение (17.31) удовлетворяло условиям (17.32), (17.33). После простых вычислений находим, что искомое решение имеет вид

$$u_n = \frac{e^{\tau\lambda} - q_2}{q_1 - q_2} q_1^n + \frac{q_1 - e^{\tau\lambda}}{q_1 - q_2} q_2^n. \quad (17.34)$$

Изучим поведение этого решения при $n \rightarrow \infty$. Пусть $t = n\tau$ фиксировано, а $\tau \rightarrow 0$. Тогда $n = t/\tau \rightarrow \infty$. С учетом (17.30) и упражнения 17.6 находим, что

$$\begin{aligned} c_1 &= \frac{e^{\tau\lambda} - q_2}{q_1 - q_2} = \frac{1 + O(\tau) + 5}{6 + O(\tau)} = 1 + O(\tau), \\ c_2 &= \frac{q_1 - e^{\tau\lambda}}{q_1 - q_2} = \frac{O(\tau^4)}{6 + O(\tau)} = O(\tau^4). \end{aligned} \quad (17.35)$$

Далее,

$$q_1^n = [e^{\tau\lambda} + O(\tau^4)]^n = e^{\lambda\tau n} (1 + O(\tau^4))^n = e^{\lambda t} (1 + O(\tau^3)). \quad (17.36)$$

Подставляя теперь (17.35), (17.36), (17.30) в (17.34), будем иметь

$$u_n = [1 + O(\tau)] e^{t\lambda} + O(\tau^4) [-5 + O(\tau)]^n.$$

Проанализируем полученный результат. Первое слагаемое аппроксимирует решение (17.26) задачи (17.25), а второе слагаемое является паразитным. Уже при не слишком больших n это слагаемое превосходит первое, ибо

$$O(\tau^4) [-5 + O(\tau)]^n = O\left(\left(\frac{t}{n}\right)^4\right) \left(-5 + O\left(\frac{t}{n}\right)\right)^n.$$

Метод (17.28) сходящимся не является.

§ 18

Устойчивость многошаговых методов

18.1 Нуль-устойчивость

Обратимся к разностному уравнению (17.27) и введем следующие обозначения

$$\rho(\zeta) := \sum_{j=0}^k \alpha_j \zeta^{k-j}, \quad \sigma(\zeta) := \sum_{j=0}^k \beta_j \zeta^{k-j}. \quad (18.1)$$

Определение 18.1. Многочлены $\rho(\zeta)$ и $\sigma(\zeta)$ из (18.1) называются соответственно первым и вторым производящими многочленами линейного многошагового метода (17.17).

Как уже было отмечено, линейный многошаговый метод (17.17) для уравнения (17.25) принимает вид линейного разностного уравнения с постоянными коэффициентами (17.27). Его характеристическое уравнение есть

$$\rho(q) - \tau \lambda \sigma(q) = 0. \quad (18.2)$$

Применительно к двушаговому методу (17.24)

$$\rho(q) = q^2 + 4q - 5,$$

а корни уравнения

$$\rho(q) = 0 \quad (18.3)$$

суть

$$q_1 = 1, \quad q_2 = -5,$$

т.е. совпадают с главными членами корней (17.30) характеристического уравнения (17.29). Именно наличие корня q_2 и привело к неустойчивости метода (17.24). Тем самым, корни уравнения (18.3) позволяют судить об устойчивости или неустойчивости метода (17.17). А они связаны с корнями характеристического уравнения (18.2). В силу (17.18), $\alpha_0 \neq 0$ и, следовательно, степени уравнений (18.2) и (18.3) совпадают.

Поэтому характеристическое уравнение (18.2) можно рассматривать как *регулярное возмущение* (при малых $\tau\lambda$) уравнения (18.3) (объяснение терминов: коэффициенты многочлена $\rho(\zeta)$ суть пределы при $\tau\lambda \rightarrow 0$ соответствующих коэффициентов характеристического многочлена, и поэтому можно говорить о возмущении; регулярность есть следствие того, что степени возмущенного и невозмущенного многочленов совпадают). Но тогда (в силу регулярности возмущения) корни уравнения (18.3) являются пределами корней уравнения (18.2) при $\tau\lambda \rightarrow 0$. Поэтому вопрос о том, будет ли решение уравнения (17.27) неограниченно возрастать при $n \rightarrow \infty$ (и фиксированном $t = n\tau$), можно решить при анализе корней уравнения (18.3). Отметим, что уравнение (18.3) является характеристическим уравнением для разностного уравнения

$$\sum_{j=0}^k \alpha_j u_{n-j} = 0, \quad (18.4)$$

которое, в свою очередь, получается из (5.27), если в нем положить $\lambda = 0$. Это означает, что (18.4) есть линейный многошаговый метод для уравнения

$$u' = 0. \quad (18.5)$$

Тем самым, отбраковка "плохих" (неустойчивых) методов может быть осуществлена при анализе их свойств применительно к уравнению (18.5).

Итак, наличие у уравнения (18.3) корней, модули которых превосходят единицу, приводит к неустойчивости. Однако опасность представляют не только такие корни, но и корни, равные по модулю единице, если они кратные. В самом деле, пусть q_1 — корень характеристического уравнения (18.3) кратности $s > 1$ такой, что $|q_1| = 1$. Тогда сеточная функция

$$P_{s-1}(n)q_1^n$$

будет растущим решением уравнения (18.4), в то время как решением уравнения (18.5), которое и аппроксимирует изучаемое уравнение (18.4), есть постоянная.

Определение 18.2. Говорят, что линейный многошаговый метод (17.17) удовлетворяет корневому условию, если

- 1) все корни первого производящего многочлена (18.1) расположены в единичном круге $|\zeta| \leq 1$;
- 2) нули $\rho(\zeta)$, расположенные на единичной окружности $|\zeta| = 1$ простые.

Определение 18.3. Линейный многошаговый метод (17.17), удовлетворяющий корневому условию, называется нуль-устойчивым (устойчивым).

Замечание 18.1. Если линейный многошаговый метод (17.17) аппроксимирует какое-либо дифференциальное уравнение, то среди нулей $\rho(\zeta)$ обязательно есть $\zeta = 1$, о чем свидетельствует первое из условий (17.20), являющее собой условие $\rho(1) = 0$.

Примеры. 1° Явный и неявный методы Адамса. В обоих случаях $\alpha_0 = 1$, $\alpha_1 = -1$, а остальные $\alpha_j = 0$. Поэтому

$$\rho(q) = q^k - q^{k-1}$$

и, следовательно,

$$q_1 = 1, \quad q_2 = \dots = q_k = 0.$$

Методы Адамса нуль-устойчивы.

2° Двухшаговая формула дифференцирования назад (17.16).

$$\begin{aligned} \rho(q) &= \frac{3}{2}q^2 - 2q + \frac{1}{2}, \\ q_1 &= 1, \quad q_2 = 1/3. \end{aligned}$$

Метод нуль-устойчив.

3° Трехшаговая формула дифференцирования назад (17.3)

$$\rho(q) = \frac{11}{6}q^3 - 3q^2 + \frac{3}{2}q - \frac{1}{3}.$$

Хотя это и многочлен третьей степени, нули его легко находятся, ибо один из его нулей есть $q_1 = 1$. Деля $\rho(q)$ на $(q - 1)$, приходим к уравнению

$$\frac{11}{6}q^2 - \frac{7}{6}q + \frac{1}{3} = 0$$

с корнями

$$q_{2,3} = \frac{7 \pm i\sqrt{39}}{22}.$$

Отсюда

$$|q_{2,3}|^2 = \frac{2}{11} < 1.$$

Метод нуль-устойчив.

Теорема 18.1 (Первый барьер Далквиста). *Порядок p устойчивого линейного k -шагового метода подчиняется следующим ограничениям:*

- $p \leq k$ для явных методов;
- $p \leq k + 1$ для неявных методов при нечетном k ;
- $p \leq k + 2$ для неявных методов при четном k .

В качестве иллюстрации первого утверждения теоремы может служить построенный нами в предыдущем параграфе явный двухшаговый метод максимального порядка аппроксимации $p = 3$, который оказался неустойчивым.

Упражнение 18.1. Построить общий явный устойчивый двухшаговый метод максимального порядка аппроксимации.

Ответ: α_0 — параметр метода,

$$\begin{aligned}\alpha_1 &= 1 - 2\alpha_0, & \alpha_2 &= \alpha_0 - 1, \\ \beta_0 &= 0, & \beta_1 &= \frac{1}{2} + \alpha_0, & \beta_2 &= \frac{1}{2} - \alpha_0.\end{aligned}$$

Условие устойчивости: $1/2 \leq \alpha_0 < \infty$. При $\alpha_0 = 1$ имеем явный метод Адамса, при $\alpha_0 = 1/2$ — метод прямоугольников с шагом $\tau' = 2\tau$. При $\alpha_0 = 1/6$ метод имеет погрешность аппроксимации $O(\tau^3)$, но неустойчив.

Упражнение 18.2. Построить устойчивый двухшаговый метод максимального порядка аппроксимации.

Ответ:

$$\begin{aligned}\alpha_0 &= 1/2, & \alpha_1 &= 0, & \alpha_2 &= -1/2, \\ \beta_0 &= 1/6, & \beta_1 &= 2/3, & \beta_2 &= 1/6.\end{aligned}$$

Этот метод иногда называется методом Симпсона (по аналогии с одноименной квадратурной формулой). Метод имеет четвертый порядок аппроксимации.

18.2 Жесткие задачи

При определении нуль-устойчивости многошагового метода мы могли ограничиться изучением простейшего дифференциального уравнения (18.5), ибо производящий многочлен $\rho(\zeta)$ из (18.1) многошагового метода (17.17), от расположения нулей которого зависит, будет ли метод устойчивым или нет, является характеристическим многочленом именно в применении к уравнению (18.5). Условие нуль-устойчивости предъявляет минимальные требования к численному методу, производя лишь грубую отбраковку абсолютно непригодных для вычислений методов. По существу, нуль-устойчивость метода обеспечивает лишь ограниченность приближенного решения для конечного временного интервала $[0, T]$ при $n \rightarrow \infty$.

Однако имеются задачи, отыскание решений которых при помощи только нуль-устойчивых методов оказывается весьма затруднительным, если не невозможным. Проще всего объяснить возникающие трудности не на примере одного уравнения, а на примере систем уравнений.

Рассмотрим однородную систему линейных дифференциальных уравнений с постоянными коэффициентами

$$\mathbf{u}' = A\mathbf{u}, \quad \mathbf{u}(0) = \mathbf{u}_0, \quad (18.6)$$

где $\mathbf{u} = [u_1 \ u_2]^T$, а

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}.$$

Найдем и проанализируем решение задачи (18.6). Как обычно, будем его искать в виде

$$\mathbf{u}(t) = \boldsymbol{\xi} e^{\lambda t}, \quad (18.7)$$

где $\boldsymbol{\xi}$ — двумерный числовой вектор, а λ — постоянная. Подставляя (18.7) в (18.6), находим, что

$$\lambda \boldsymbol{\xi} e^{\lambda t} = e^{\lambda t} A \boldsymbol{\xi},$$

а, сокращая на $e^{\lambda t}$, получим следующую задачу на собственные значения:

$$A \boldsymbol{\xi} = \lambda \boldsymbol{\xi}. \quad (18.8)$$

Будем предполагать, что A — матрица простой структуры, т.е. у нее имеется полный набор собственных векторов. Тогда

$$A \boldsymbol{\xi}_1 = \lambda_1 \boldsymbol{\xi}_1, \quad A \boldsymbol{\xi}_2 = \lambda_2 \boldsymbol{\xi}_2$$

и $\boldsymbol{\xi}_1$ и $\boldsymbol{\xi}_2$ линейно независимы.

В рассматриваемом случае общее решение системы (18.6) принимает вид

$$\mathbf{u}(t) = c_1 \boldsymbol{\xi}_1 e^{\lambda_1 t} + c_2 \boldsymbol{\xi}_2 e^{\lambda_2 t}, \quad (18.9)$$

а решение задачи Коши (18.6) получается отсюда при значениях c_1 и c_2 , найденных из алгебраической системы

$$\boldsymbol{\xi}_1 c_1 + \boldsymbol{\xi}_2 c_2 = \mathbf{u}_0. \quad (18.10)$$

Будем для простоты предполагать, что собственные числа λ_1 и λ_2 действительны. Более существенным для нас будет предположение об их отрицательности

$$\lambda_1 < 0, \quad \lambda_2 < 0. \quad (18.11)$$

В силу сделанных предположений модули компонент u_1 и u_2 решения (18.9) будут стремиться к нулю при $t \rightarrow \infty$.

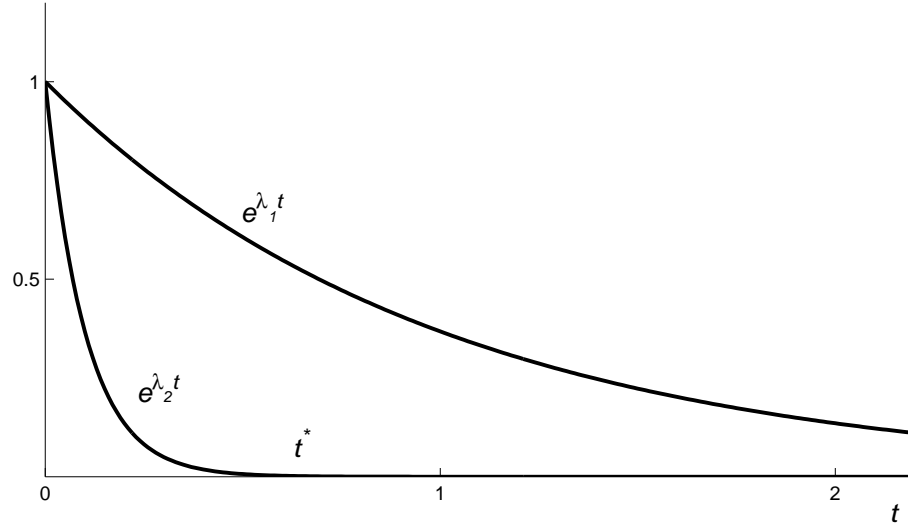


Рис. 1

Предположим теперь дополнительно, что

$$\lambda_1 = O(1), \quad |\lambda_2| \gg |\lambda_1|. \quad (18.12)$$

Так как в этом случае $e^{\lambda_2 t}$ убывает значительно быстрее $e^{\lambda_1 t}$, то через некоторое время t^* составляющая $c_2 \xi_2 e^{\lambda_2 t}$ решения (18.9) будет практически равной нулю, и решение будет почти полностью определяться составляющей $c_1 \xi_1 e^{\lambda_1 t}$. (см. рис. 1)

В рассматриваемой ситуации естественно было бы ожидать, что и у численного решения задачи (18.6) модули компонент хотя бы не возрастали.

Применим для решения задачи (18.6) метод Эйлера

$$\frac{\mathbf{u}^{n+1} - \mathbf{u}^n}{\tau} = A\mathbf{u}^n, \quad \mathbf{u}^0 = \mathbf{u}_0. \quad (18.13)$$

Найдем решение задачи (18.13). Искать его будем в виде (см. (6.30))

$$\mathbf{u}^n = \boldsymbol{\xi} q^n, \quad q = \text{const} \neq 0. \quad (18.14)$$

Подставляя (18.14) в (18.13), получим

$$q^n \frac{q - 1}{\tau} \boldsymbol{\xi} = q^n A \boldsymbol{\xi},$$

а после сокращения на q^n обнаруживаем, что для отыскания ξ имеем задачу (18.8) с $\lambda = (q - 1)/\tau$. Поэтому $q = 1 + \tau\lambda$, и решение задачи (18.13) есть

$$\mathbf{u}^n = c_1 \xi_1 (1 + \tau\lambda_1)^n + c_2 \xi_2 (1 + \tau\lambda_2)^n, \quad (18.15)$$

где c_1, c_2 — решение системы (18.10).

Чтобы модули компонент решения (18.15) не возрастали при $n \rightarrow \infty$, необходимо и достаточно, чтобы выполнялись условия

$$|1 + \tau\lambda_1| \leq 1, \quad |1 + \tau\lambda_2| \leq 1,$$

что вместе с (18.11) и (18.12) приводит к условию

$$\tau \leq 2/|\lambda_2| \ll 1. \quad (18.16)$$

Ограничение (18.16), вообще говоря, является довольно жестким. Если при $t \leq t^*$ это ограничение вполне разумно, и даже из соображений аппроксимации и точности нужно требовать $\tau \ll 2/|\lambda_2|$, то при $t > t^*$, когда вторая составляющая каждой компоненты решения (18.15) вроде бы не должна поставлять новой информации, и желательно было бы увеличить шаг τ с той целью, чтобы сэкономить ресурсы и не воспроизводить первую составляющую с излишней точностью. Но тогда придется нарушить условие (18.16), что приведет к резкому возрастанию второй составляющей решения и полной потере точности.

Определение 18.4. Система дифференциальных уравнений (18.6) с постоянной матрицей A порядка m называется жесткой, если

1° $\operatorname{Re} \lambda_j < 0$, $j = 0, \dots, m$,

2° отношение

$$S = \frac{\max_j |\operatorname{Re} \lambda_j|}{\min_j |\operatorname{Re} \lambda_j|} \gg 1. \quad (18.17)$$

Определение 18.5. Число S из (18.17) называется коэффициентом жесткости задачи (18.6).

Замечание 18.2. Для линейной системы с матрицей A , зависящей от t , коэффициент жесткости также зависит от t , и, если он велик для каких-либо t из интересующего нас интервала, то система жесткая. Для нелинейных систем жесткость определяется в окрестности какого-либо решения при помощи соответствующей матрицы Якоби.

Применим теперь для решения задачи (18.6) неявный метод Эйлера

$$\frac{\mathbf{u}^{n+1} - \mathbf{u}^n}{\tau} = A \mathbf{u}^{n+1}.$$

Подставляя сюда (18.14), находим, что

$$q^{n+1} \frac{1 - q^{-1}}{\tau} \xi = q^{n+1} A \xi,$$

т.е. $\lambda\tau = (1 - q^{-1})$, $q = (1 - \tau\lambda)^{-1}$ и

$$\mathbf{u}^n = c_1 \boldsymbol{\xi}_1 (1 - \tau\lambda)^{-n} + c_2 \boldsymbol{\xi}_2 (1 - \tau\lambda_2)^{-n}.$$

Очевидно, что при выполнении условий (18.11) модули компонент \mathbf{u}^n монотонно убывают при $n \rightarrow \infty$ при *любых* τ , и, следовательно, τ можно выбирать только из соображений точности.

Неявный метод Эйлера при решении жестких систем оказался существенно более устойчивым, чем просто метод Эйлера.

Как отобрать методы, пригодные для решения жестких задач? Ужесточить требование устойчивости.

18.3 A-устойчивость

Если при определении нуль-устойчивости основной моделью было уравнение (18.5), то теперь следует обратиться к уравнению (17.25). Многошаговый метод (17.17) в применении к линейному однородному уравнению (17.25) имеет вид (17.27), а характеристическое уравнение этого разностного уравнения задается соотношением (18.2).

Определение 18.6. Линейный многошаговый метод (17.17) в применении к уравнению (17.25) называется абсолютно устойчивым для данного λ и данного τ , если при указанном значении $\tau\lambda$ все корни характеристического уравнения (18.2) расположены внутри единичного круга.

Определение 18.7. Множество всех точек комплексной плоскости $\tau\lambda$, для которых линейный многошаговый метод (17.17) в применении к (17.25) абсолютно устойчив, называется областью абсолютной устойчивости метода.

Пример 4°. Метод Эйлера (15.7). Единственный корень характеристического уравнения $q = 1 + \tau\lambda$. Условие абсолютной устойчивости

$$|1 + \tau\lambda| \leq 1.$$

Областью абсолютной устойчивости является единичный круг с центром в точке $\tau\lambda = -1$. (см. рис. 2)

Пример 5°. Неявный метод Эйлера (15.8). Условие абсолютной устойчивости

$$|q| = |1 - \tau\lambda|^{-1} \leq 1, \quad \text{т.е.} \quad |1 - \tau\lambda| \geq 1.$$

Областью абсолютной устойчивости является внешность единичного круга с центром в точке $\tau\lambda = 1$. (см. рис. 3)

Определение 18.8. Линейный многошаговый метод (17.17) называется A-устойчивым, если область его абсолютной устойчивости содержит левую полуплоскость $\text{Re}(\tau\lambda) < 0$.

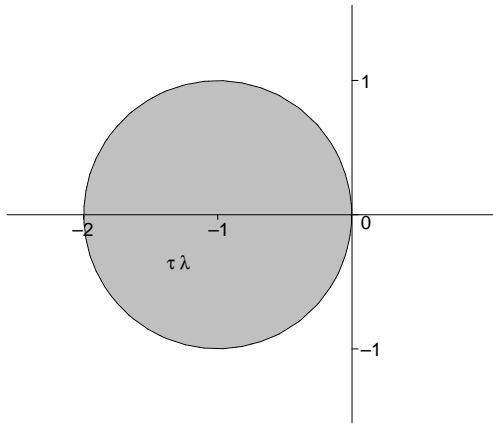


Рис. 2

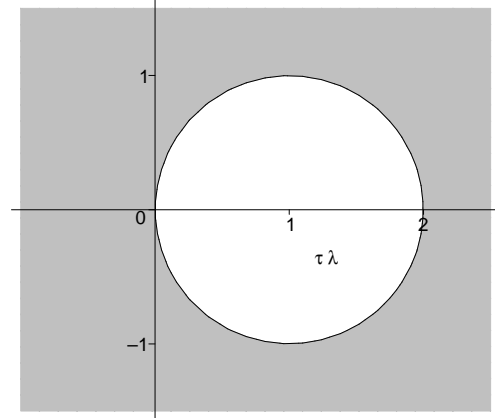


Рис. 3

Из приведенных примеров следует, что метод Эйлера не является A -устойчивым, а неявный метод Эйлера A -устойчив.

Пример 6°. Метод трапеций (15.12). Применительно к уравнению (17.25) этот метод имеет вид

$$\frac{u^{n+1} - u^n}{\tau} = \lambda \frac{u^{n+1} + u^n}{2},$$

а его характеристическое уравнение есть $(q - 1)/\tau = \lambda(q + 1)/2$. Отсюда находим единственный корень

$$q = \frac{1 + \tau\lambda/2}{1 - \tau\lambda/2}$$

и условие абсолютной устойчивости

$$|q| = \left| \frac{1 + \tau\lambda/2}{1 - \tau\lambda/2} \right| \leq 1$$

или

$$|1 + \tau\lambda/2| \leq |1 - \tau\lambda/2|.$$

Пусть $\tau\lambda = x + iy$. Тогда условие абсолютной устойчивости примет вид

$$\left| 1 + \frac{x}{2} + i\frac{y}{2} \right| \leq \left| 1 - \frac{x}{2} - i\frac{y}{2} \right|.$$

или

$$\left(1 + \frac{x}{2} \right)^2 + \frac{y^2}{4} \leq \left(1 - \frac{x}{2} \right)^2 + \frac{y^2}{4}.$$

Раскрывая скобки, находим, что условие абсолютной устойчивости есть

$$x = \operatorname{Re}(\tau\lambda) < 0.$$

Областью абсолютной устойчивости метода трапеций является левая полуплоскость $\operatorname{Re}(\tau\lambda) < 0$ (Рис. 4). Метод A -устойчив.

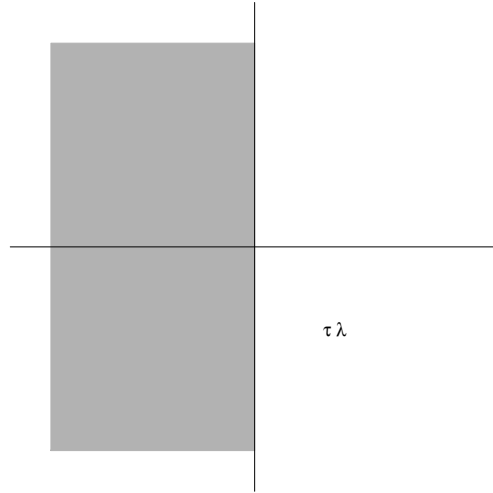


Рис. 4

Теорема 18.2. Среди линейных многошаговых методов (17.17) не существует явных A -устойчивых методов.

Теорема 18.3. Среди неявных линейных многошаговых методов (17.17) не существует A -устойчивых методов, имеющих порядок точности выше второго.

Пример 7°. Двухшаговая формула дифференцирования назад. Этот метод задается соотношением (17.16)

$$\left(\frac{3}{2}u_{n+1} - 2u_n + \frac{1}{2}u_{n-1}\right) = \tau f(u_{n+1}). \quad (18.18)$$

Характеристическое уравнение, отвечающее этому методу в применении к уравнению (17.25) есть

$$\frac{3}{2}q^2 - 2q + \frac{1}{2} - \tau\lambda q^2 = 0. \quad (18.19)$$

Определим область абсолютной устойчивости этого метода. Для этого достаточно найти ее границу, т.е. такое множество комплексной плоскости $z = \tau\lambda$, где $|q(z)| = 1$. С этой целью выразим из (18.19) $\tau\lambda$ через q

$$z = \frac{3}{2} - \frac{2}{q} + \frac{1}{2q^2}. \quad (18.20)$$

Поскольку нас интересуют значения $|q| = 1$, то пусть $q = e^{-i\varphi}$. Отсюда и из (18.20)

$$z = \frac{3}{2} - 2e^{i\varphi} + \frac{1}{2}e^{2i\varphi}. \quad (18.21)$$

При изменении аргумента φ от 0 до 2π точка z из (18.21) описывает замкнутую кривую, симметричную относительно действительной оси (функция $\sin k\varphi$ — нечетная), которая и является границей области абсолютной устойчивости.

$$\begin{aligned} z &= \frac{3}{2} - 2 \cos \varphi + \frac{1}{2} \cos 2\varphi + i(-2 \sin \varphi + \frac{1}{2} \sin 2\varphi) = \\ &= \frac{3}{2} - 2 \cos \varphi + \cos^2 \varphi - \frac{1}{2} + i(-2 \sin \varphi + \sin \varphi \cos \varphi) = \\ &= (1 - \cos \varphi)^2 \pm i\sqrt{1 - \cos^2 \varphi}(2 - \cos \varphi) = \\ &= (1 - t)^2 \pm i\sqrt{1 - t^2}(2 - t), \quad t = \cos \varphi. \end{aligned}$$

Отсюда следует, что

$$\operatorname{Re} z = (1 - t)^2 \geq 0,$$

и, следовательно, кривая расположена в правой полуплоскости. Построим ее. Мнимая часть $z(t)$ обращается в нуль при $t = \pm 1$. Действительная часть $z(t)$ при этих значениях параметра равна 0 и 4.

Исследования показывают, что

$$\begin{aligned} \max_{[-1,1]} \operatorname{Im} z(t) &= \operatorname{Im} z \left(\frac{1 - \sqrt{3}}{2} \right) = \frac{(3 + \sqrt{3})\sqrt[4]{3}}{2\sqrt{2}} \approx 2.20, \\ \operatorname{Re} z \left(\frac{1 - \sqrt{3}}{2} \right) &= \frac{2 + \sqrt{3}}{2} \approx 1.86. \end{aligned}$$

Из (18.20) находим, что при

$$|q| \rightarrow \infty, \quad z \rightarrow \frac{3}{2} \in G,$$

и, следовательно, внутренность области — область неустойчивости. Тем самым, вне G (Рис. 5) $|q| < 1$, и метод абсолютно устойчив, а, следовательно, и А-устойчив. Этот метод второго порядка точности.

Пример 8°. Трехшаговая формула дифференцирования назад. (Упражнение 17.3)

$$\frac{11}{6}u_{n+1} - 3u_n + \frac{3}{2}u_{n-1} - \frac{1}{3}u_{n-2} = \tau\lambda u_{n+1}. \quad (18.22)$$

Характеристическое уравнение этого разностного уравнения имеет вид

$$\frac{11}{6}q^3 - 3q^2 + \frac{3}{2}q - \frac{1}{3} = \tau\lambda q^3. \quad (18.23)$$

Снова положим $|q| = 1$, т.е. $q = e^{-i\varphi}$ и $\tau\lambda = z$. Тогда

$$z = \frac{11}{6} - 3e^{i\varphi} + \frac{3}{2}e^{2i\varphi} - \frac{1}{3}e^{3i\varphi}.$$

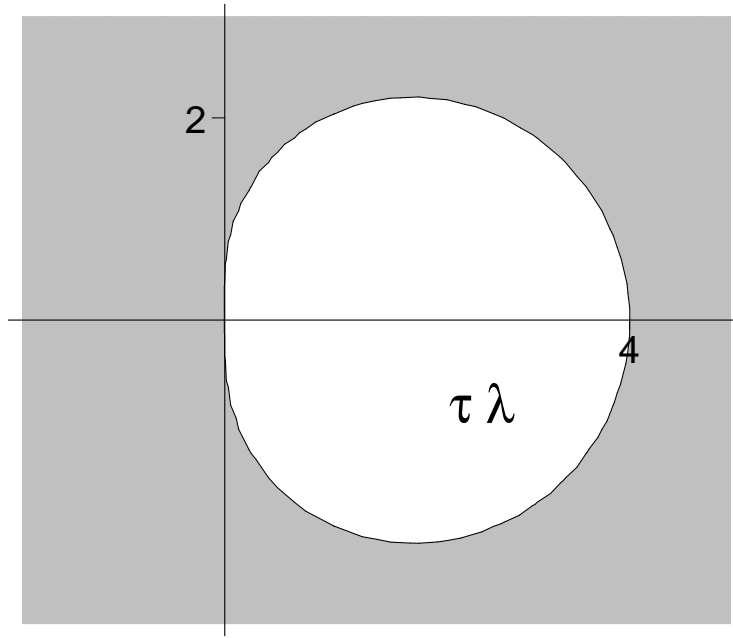


Рис. 5

Обозначая $\cos \varphi = t$, после простых вычислений находим, что

$$z = -\frac{1}{3}(t-1)^2(4t-1) \pm \frac{i}{3}\sqrt{1-t^2}(4t^2-9t+8).$$

При $t = \pm 1$ $\operatorname{Im} z = 0$, а $\operatorname{Re} z = 0$ или $20/3$. Исследования показывают, что $\operatorname{Re} z$ как функция t принимает экстремальные значения при $t = 1/2$ и $t = 1$. Значение $t = 1$ мы уже рассмотрели, а

$$\operatorname{Re} z(1/2) = \min \operatorname{Re}(t) = -1/12, \quad \operatorname{Im} z(1/2) = \pm 3\sqrt{3}/4 \approx \pm 1.30$$

и, следовательно, часть границы области устойчивости расположена в левой полуплоскости. Как легко видеть, мнимую ось граница устойчивости пересекает при $t = 1/4$ и

$$\operatorname{Im} z(1/4) = \pm\sqrt{15}/2 \approx \pm 1.94.$$

Экстремальные значения $\operatorname{Im} z(t)$ принимает в точке

$$t^* = -\frac{1}{2} \left[(2 + \sqrt{3})^{1/3} + (2 + \sqrt{3})^{-1/3} - 1 \right] \approx -0.60,$$

причем

$$\operatorname{Im} z(t^*) \approx \pm 3.96, \quad \operatorname{Re} z(t^*) \approx 2.89.$$

Область устойчивости изображена на рис. 6.

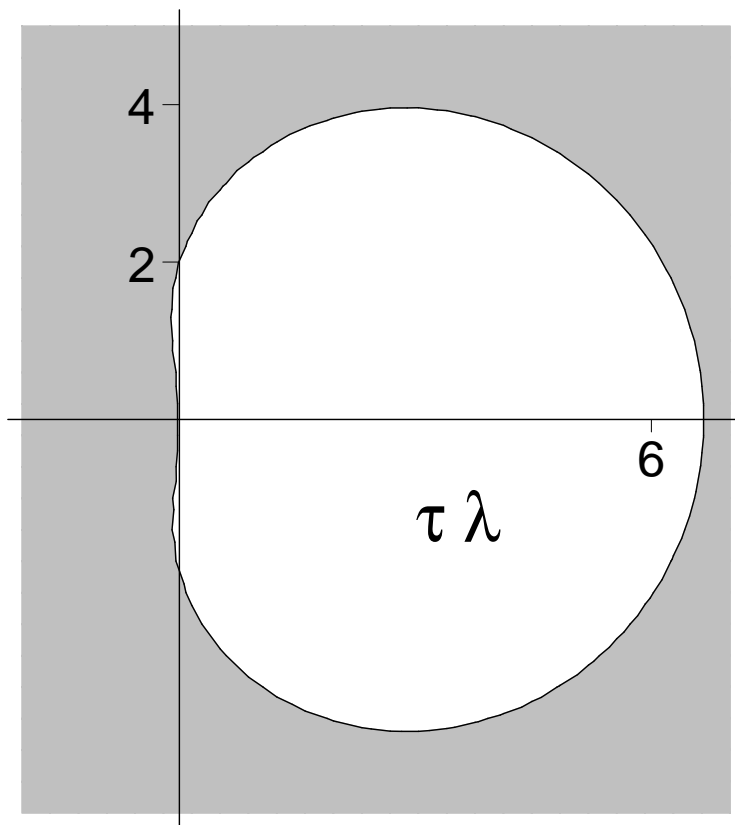


Рис. 6.

Определение 18.9. Линейный многошаговый метод называется $A(\alpha)$ -устойчивым, если его область абсолютной устойчивости содержит угол

$$|\arg(-\tau\lambda)| < \alpha.$$

Замечание 18.3. $A(\pi/2)$ - и A -устойчивости совпадают.

Теорема 18.4. Существуют многошаговые методы 3-го и 4-го порядков точности $A(\alpha)$ -устойчивые при любых $\alpha < \pi/2$.

Теорема 18.5. Явные линейные многошаговые методы не являются $A(\alpha)$ -устойчивыми ни при каких α .

Теорема 18.6. Методы дифференцирования назад при $k \leq 6$ являются $A(\alpha)$ -устойчивыми при соответствующих значениях $\alpha \neq 0$.

Упражнение 18.3. Исследовать область абсолютной устойчивости двухшагового неявного метода Адамса.

18.4 Устойчивость методов Рунге-Кутты

Как было уже отмечено, методы (все) Рунге-Кутты являются нуль-устойчивыми. Исследуем области абсолютной устойчивости некоторых из этих методов. Рассмотрим явный трехэтапный метод третьего порядка, задаваемый таблицей (16.44), которая имеет вид

$$\begin{array}{c|ccc} 1/2 & 1/2 & & \\ 1 & -1 & 2 & \\ \hline & 1/6 & 2/3 & 1/6 \end{array}$$

Применительно к уравнению (17.25) этот метод задается следующими соотношениями

$$\begin{aligned} Y_1 &= u_n, \\ Y_2 &= u_n + \frac{\tau\lambda}{2}Y_1, \\ Y_3 &= u_n - \tau\lambda Y_1 + 2\tau\lambda Y_2, \\ u_{n+1} &= u_n + \tau\lambda \left(\frac{1}{6}Y_1 + \frac{2}{3}Y_2 + \frac{1}{6}Y_3 \right). \end{aligned}$$

Исключая из этих соотношений промежуточные величины Y_1 , Y_2 и Y_3 , будем иметь

$$\begin{aligned} Y_2 &= \left(1 + \frac{\tau\lambda}{2} \right) u_n, \\ Y_3 &= \left[1 - \tau\lambda + 2\tau\lambda \left(1 + \frac{\tau\lambda}{2} \right) \right] u_n = (1 + \tau\lambda + \tau^2\lambda^2)u_n, \\ u_{n+1} &= \left\{ 1 + \tau\lambda \left[\frac{1}{6} + \frac{2}{3} \left(1 + \frac{\tau\lambda}{2} \right) + \frac{1}{6}(1 + \tau\lambda + \tau^2\lambda^2) \right] \right\} u_n = \\ &= \left(1 + \tau\lambda + \frac{\tau^2\lambda^2}{2} + \frac{\tau^3\lambda^3}{6} \right) u_n. \end{aligned}$$

Это есть линейное разностное уравнение первого порядка, единственный корень характеристического уравнения которого равен

$$q = 1 + \tau\lambda + \frac{\tau^2\lambda^2}{2} + \frac{\tau^3\lambda^3}{6} = e^{\tau\lambda} + O(\tau^4\lambda^4).$$

Обозначим $\tau\lambda$ через z . Тогда

$$q = 1 + z + \frac{z^2}{2} + \frac{z^3}{6}.$$

Этот корень характеристического уравнения есть многочлен третьей степени от z и в левой полуплоскости $\operatorname{Re} z < 0$ ограниченным быть не может. Метод не является $A(\alpha)$ -устойчивым ни при каком α .

Рассмотрим теперь двухэтапный метод третьего порядка, задаваемый таблицей (16.33), которая имеет вид

$$\frac{\begin{array}{c|cc} \theta_1 = \gamma & \gamma & 0 \\ \theta_2 = 1 - \gamma & 1 - 2\gamma & \gamma \end{array}}{\begin{array}{c|cc} & 1/2 & 1/2 \end{array}} \quad \gamma = \frac{3 \pm \sqrt{3}}{6}.$$

Применительно к уравнению (17.25) этот метод записывается следующим образом

$$\begin{aligned} Y_1 &= u_n + \gamma\tau\lambda Y_1, \\ Y_2 &= u_n + \tau\lambda(1 - 2\gamma)Y_1 + \tau\lambda\gamma Y_2, \\ u_{n+1} &= u_n + \frac{\tau\lambda}{2}(Y_1 + Y_2). \end{aligned}$$

Как и в предыдущем примере, положим $\tau\lambda = z$ и исключим Y_1 и Y_2 . Решая систему линейных алгебраических уравнений второго порядка относительно Y_1 и Y_2 (первые два уравнения) и подставляя результат в третье уравнение, находим, что

$$\begin{aligned} Y_1 &= \frac{u_n}{1 - \gamma z}, \quad Y_2 = \frac{1 + (1 - 3\gamma)z}{(1 - \gamma z)^2} u_n, \\ u_{n+1} &= \left[1 + \frac{z}{2} \left(\frac{1}{1 - \gamma z} + \frac{1 + (1 - 3\gamma)z}{(1 - \gamma z)^2} \right) \right] u_n. \end{aligned}$$

Отсюда следует, что единственным корнем характеристического уравнения является

$$\begin{aligned} q &= \frac{1 - 2\gamma z + \gamma^2 z^2 + z/2 - \gamma z^2/2 + z/2 + (1 - 3\gamma)z^2/2}{(1 - \gamma z)^2} = \\ &= \frac{1 + (1 - 2\gamma)z + (\gamma^2 - 2\gamma + 1/2)z^2}{1 - 2\gamma z + \gamma^2 z^2} = \frac{P(z)}{Q(z)}. \end{aligned}$$

Этот корень является дробно-рациональной функцией, полюсом второго порядка которой является точка $z = \gamma^{-1} = (3 \mp \sqrt{3})/6$, расположенная в правой полуплоскости. В левой полуплоскости эта функция аналитична и, следовательно, максимум ее модуля здесь не превосходит максимума модуля на границе, т.е. при $z = iy$. Оценим ее модуль на мнимой оси. Имеем

$$\begin{aligned} |P(iy)|^2 &= [1 - (\gamma^2 - 2\gamma + 1/2)y^2]^2 + (1 - 2\gamma)^2 y^2 = \\ &= 1 - 2(\gamma^2 - 2\gamma + 1/2)y^2 + (\gamma^2 - 2\gamma + 1/2)^2 y^4 + (1 - 4\gamma + 4\gamma^2)y^2 = \\ &= 1 + 2\gamma^2 y^2 + (\gamma^2 - 2\gamma + 1/2)^2 y^4 \end{aligned}$$

и

$$|Q(iy)|^2 = (1 - \gamma^2 y^2)^2 + 4\gamma^2 y^2 = 1 + 2\gamma^2 y^2 + \gamma^4 y^4.$$

Отсюда

$$\left| \frac{P(iy)}{Q(iy)} \right|^2 = \frac{1 + 2\gamma^2 y^2 + (\gamma^2 - 2\gamma + 1/2)^2 y^4}{1 + 2\gamma^2 y^2 + \gamma^4 y^4}.$$

Добавим к числителю и вычтем из него $\gamma^4 y^4$, после чего выделим единицу

$$\left| \frac{P(iy)}{Q(iy)} \right|^2 = 1 + \frac{(\gamma^2 - 2\gamma + 1/2 - \gamma^2)(\gamma^2 - 2\gamma + 1/2 + \gamma^2)y^4}{(1 + \gamma^2 y^2)^2}.$$

Подставим вместо γ его значения из (16.33). Найдем, что

$$-2\gamma + 1/2 = \frac{-3 \mp 2\sqrt{3}}{6},$$

а

$$2\gamma^2 - 2\gamma + 1/2 = \frac{9 + 3 \pm 6\sqrt{3}}{18} + \frac{-3 \mp 2\sqrt{3}}{6} = \frac{1}{6}.$$

Поэтому

$$\left| \frac{P(iy)}{Q(iy)} \right|^2 = 1 - \frac{3 \pm 2\sqrt{3}}{36} \frac{y^4}{(1 + \gamma^2 y^2)^2}.$$

Поскольку это выражение не меньше нуля, а при $\gamma = (3 + \sqrt{3})/6$ (верхний знак в коэффициенте у второго слагаемого) вычитаемое неотрицательно, то в рассматриваемом случае

$$\left| \frac{P(iy)}{Q(iy)} \right| \leq 1,$$

и изучаемый метод является A -устойчивым. При $\gamma = (3 - \sqrt{3})/6$ вычитаемое отрицательно, и поэтому $|P(iy)/Q(iy)| > 1$ для $y \neq 0$. В этом случае метод A -устойчивым не является. Области абсолютной устойчивости этих методов изображены на рисунках 7 и 8, соответственно.

Тем самым, один из методов (16.33), именно, отвечающий $\gamma = (3 + \sqrt{3})/6$, является A -устойчивым, в то время как второй таким свойством не обладает и даже не является $A(\alpha)$ -устойчивым.

Упражнение 18.4. Доказать, что неявный двухэтапный метод Рунге-Кутты четвертого порядка (оптимальный двухэтапный метод) (16.34) является $A(\alpha)$ -устойчивым.

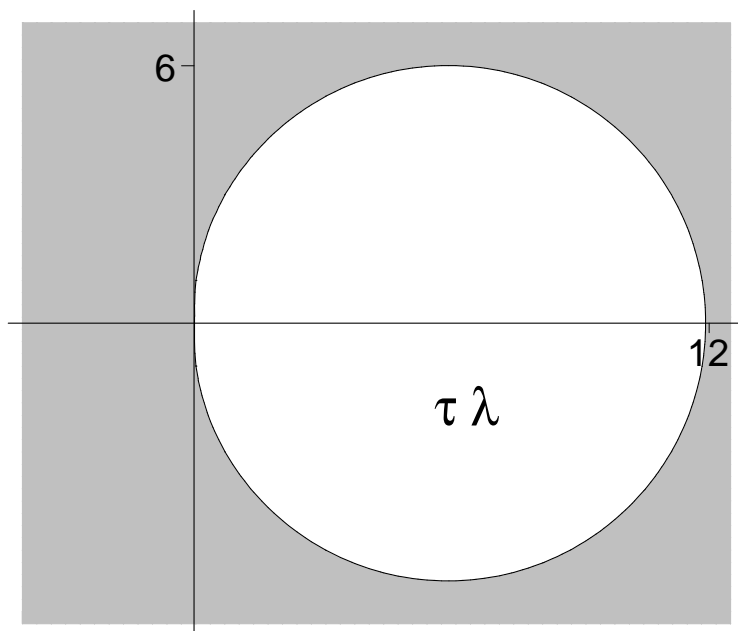


Рис. 7

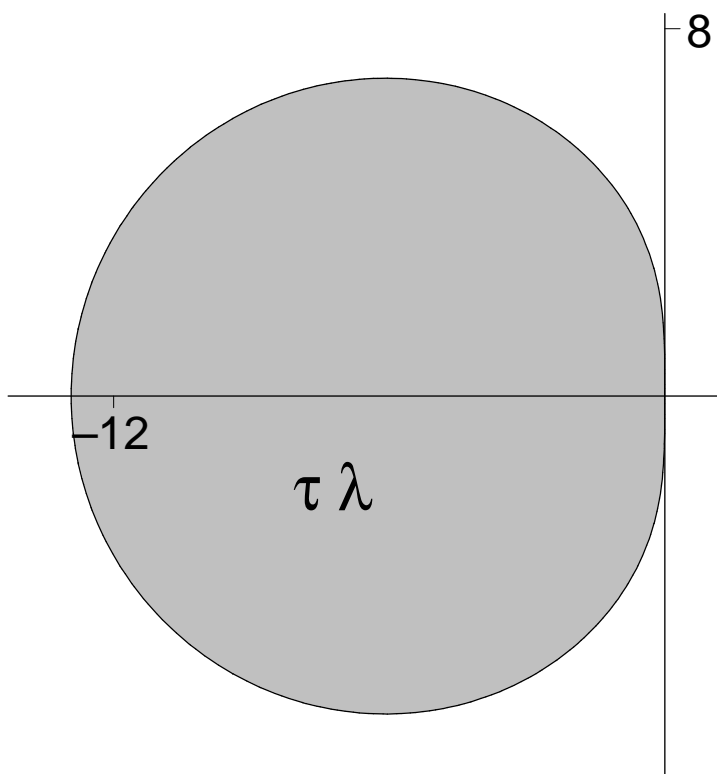


Рис. 8

Глава V

Численные методы решения краевых задач для обыкновенных дифференциальных уравнений

§ 19

Элементы теории разностных схем

19.1 Введение

Простейшим содержательным примером краевой задачи для обыкновенного дифференциального уравнения является следующий

$$-u''(x) = f(x), \quad 0 < x < l, \quad (19.1)$$

$$u(0) = g_0, \quad u(l) = g_1. \quad (19.2)$$

У краевой задачи, в отличие от задачи Коши, дополнительные условия, выделяющие единственное решение уравнения (19.1), задаются не в одной точке, а в нескольких (обычно в двух), и называются краевыми (или граничными) условиями. Это вносит дополнительные трудности в процесс решения задачи.

Мы будем изучать разностные методы решения краевых задач. Для этого на отрезке $[0, l]$ введем сетку

$$\bar{\omega} := \{x = x_i = ih \mid i = 0, \dots, N\}.$$

Точки x_i будем называть узлами сетки, а число $h = l/N$ — ее шагом. Введенная сетка является равномерной. Если бы расстояния между узлами менялись при переходе от одного узла к другому, то сетка была бы неравномерной.

Суть разностных методов решения краевых задач для дифференциальных уравнений состоит в том, что производные, входящие в дифференциальное уравнение и граничные условия, заменяются подходящими разностными отношениями. В результате краевая задача заменяется (аппроксимируется) системой алгебраических (линейных, если исходная задача была линейной) уравнений, решение которой и принимается за приближенное решение краевой задачи.

Напомним простейшие аппроксимации первой и второй производных

$$\frac{u(x_i) - u(x_{i-1}))}{h} = u'(x_i) + O(h), \quad (19.3)$$

$$\frac{u(x_{i+1}) - u(x_i)}{h} = u'(x_i) + O(h), \quad (19.4)$$

$$\frac{u(x_{i+1}) - u(x_{i-1}))}{2h} = u'(x_i) + O(h^2), \quad (19.5)$$

$$\frac{-u(x_{i+2}) + 4u(x_{i+1}) - 3u(x_i))}{2h} = u'(x_i) + O(h^2), \quad (19.6)$$

$$\frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))}{h^2} = u''(x_i) + O(h^2). \quad (19.7)$$

Для справедливости соотношений (19.3) и (19.4) достаточно, чтобы $u(x) \in C^2$, для справедливости (19.5) и (19.6) — $u(x) \in C^3$, для справедливости (19.7) — $u(x) \in C^4$. В этом можно убедиться путем разложения левых частей (19.3)-(19.7) в точке $x = x_i$ по формуле Тейлора.

Упражнение 19.1. Убедиться в справедливости (19.3)-(19.7).

Замечание 19.1. Если функцию $u(x)$ заменить интерполяционным многочленом Лагранжа первой степени по узлам x_{i-1} и x_i или x_i и x_{i+1} , а затем его продифференцировать, то получим левые части соотношений (19.3), (19.4). Заменяя $u(x)$ интерполяционным многочленом второй степени по узлам x_{i-1}, x_i, x_{i+1} или x_i, x_{i+1}, x_{i+2} , дифференцируя полученный интерполянт и полагая $x = x_i$, получим левые части (19.5) и (19.6), соответственно.

Воспользуемся соотношением (19.7) для замены второй производной в (19.1) разностным отношением

$$-\frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))}{h^2} \approx f(x_i), \quad x_i = h, 2h, \dots, l-h.$$

Превратим приближенные равенства в точные путем замены точного решения $u(x_i)$ в узле x_i на приближенное u_i^h :

$$-\frac{u_{i+1}^h - 2u_i^h + u_{i-1}^h}{h^2} = f_i, \quad i = \overline{1, N-1}, \quad (19.8)$$

Это есть система $(N-1)$ линейных алгебраических уравнений с $N+1$ неизвестными $u_0^h, u_1^h, \dots, u_N^h$. Система (19.8) недоопределена (как и следовало ожидать). Воспользуемся граничными условиями (19.2) и положим

$$u_0^h = g_0, \quad u_N^h = g_1. \quad (19.9)$$

Решение системы (19.8), (19.9), если оно существует, будем называть приближенным решением задачи (19.1), (19.2).

19.2 Основные понятия теории разностных схем

Обозначим дифференциальное выражение, стоящее в левой части (19.1), через Lu . Тогда дифференциальное уравнение (19.1) примет вид

$$Lu = f(x), \quad 0 < x < l. \quad (19.10)$$

Граничные условия (19.2) запишем в виде

$$lu = g. \quad (19.11)$$

Аналогично, разностное выражение, стоящее в левой части (19.8), обозначим через $L^h u^h$. Тогда из (19.8) будем иметь

$$L^h u_i^h = f_i^h, \quad i = \overline{1, N-1}, \quad (19.12)$$

где $f_i^h = f_i$. Граничные условия (19.9) запишем в виде, аналогичном (19.11)

$$l^h u^h = g^h. \quad (19.13)$$

Определение 19.1. Сеточная функция

$$\Psi_v(x) := L^h v - Lv, \quad x \in \omega, \quad (19.14)$$

определенная на сетке ω , где v — достаточно гладкая функция, заданная на $[0, l]$, называется погрешностью аппроксимации дифференциального выражения Lv разностным выражением $L^h v$.

Определение 19.2. Разностное выражение $L^h v$ аппроксимирует дифференциальное выражение Lv , если погрешность аппроксимации $\Psi_v \rightarrow 0$ (в каком-нибудь смысле) при $h \rightarrow 0$.

Определение 19.3. Сеточная функция

$$z = u^h - u, \quad x \in \bar{\omega}, \quad (19.15)$$

где u^h — решение задачи (19.12), (19.13), а u — решение задачи (19.10), (19.11), называется погрешностью решения.

Сформулируем задачу для погрешности решения z . Подставим в (19.12), (19.13) u^h , выражаемое из (19.15) через z и u : $u^h = z + u$. Будем иметь

$$L^h z = f^h - L^h u, \quad l^h z = g^h - l^h u. \quad (19.16)$$

Определение 19.4. Функция

$$\Psi = f^h - L^h u, \quad x \in \omega, \quad (19.17)$$

являющаяся правой частью уравнения для погрешности решения (19.16), называется погрешностью аппроксимации уравнения (19.10) уравнением (19.12).

Определение 19.5. Функция

$$\psi = g^h - l^h u, \quad (19.18)$$

являющаяся правой частью в граничных условиях для погрешности решения (19.16), называется погрешностью граничных условий (19.11) граничными условиями (19.13).

Замечание 19.2. Так как в силу (19.10) $Lu - f = 0$, то, добавляя этот нуль к представлению погрешности аппроксимации (19.17), будем иметь

$$\Psi = f^h - L^h u = f^h - f - (L^h u - Lu) = (f^h - f) - \Psi_u, \quad (19.19)$$

где Ψ_u определяется соотношением (19.14). Тем самым, погрешность аппроксимации уравнения представляет собой разность между погрешностью аппроксимации правой части и погрешностью аппроксимации дифференциального выражения. Аналогичные представления имеют место и для погрешности аппроксимации граничных условий:

$$\psi = g^h - l^h u = g^h - g - (l^h u - lu) = (g^h - g) - \psi_u. \quad (19.20)$$

Определение 19.6. Задача (19.12), (19.13) аппроксимирует задачу (19.10), (19.11), если Ψ и ψ стремятся к нулю при $h \rightarrow 0$ вместе с Ψ_u и ψ_u .

Определение 19.7. Решение задачи (19.12), (19.13) сходится к решению задачи (19.10), (19.11), если $z \rightarrow 0$ (в каком-либо смысле) при $h \rightarrow 0$.

Определение 19.8. Задача (19.12), (19.13) аппроксимирует задачу (19.10), (19.11) с погрешностью порядка $n > 0$, если

$$\|\Psi_u\|_{(1)} = o(1), \quad \|\psi_u\|_{(2)} = o(1), \quad \|\Psi\|_{(1)} = O(h^n), \quad \|\psi\|_{(2)} = O(h^n)$$

Определение 19.9. Решение задачи (19.12), (19.13) сходится к решению задачи (19.10), (19.11) со скоростью $O(h^n)$, если

$$\|z\|_{(3)} = O(h^n).$$

Проиллюстрируем введенные понятия на примере задачи (19.1), (19.2). Так как в данном случае $L = -d^2v/dx^2$, а

$$L^h v = -\frac{v(x_{i+1}) - 2v(x_i) + v(x_{i-1}))}{h^2},$$

то, в силу (19.7),

$$\Psi_v = O(h^2),$$

т.е. дифференциальное выражение v'' аппроксимируется разностным выражением $(v_{i+1} - 2v_i + v_{i-1})/h^2$ на функциях $v(x) \in C^4$ с погрешностью $O(h^2)$.

Далее, так как $f_i^h = f(x_i)$, то с учетом (19.19) заключаем, что дифференциальное уравнение (19.1) аппроксимируется разностным уравнением (19.8) с погрешностью $O(h^2)$, если $u(x) \in C^4[0, l]$.

Наконец,

$$\begin{aligned} lu &= \{u(0), u(1)\}, \\ l^h u &= \{u_0, u_N\}, \\ g &= \{g_0, g_1\} = g^h, \end{aligned}$$

так что

$$\psi = g^h - l^h u = 0.$$

Итак, задача (19.8), (19.9) аппроксимирует задачу (19.1), (19.2) (при $u(x) \in C^4[0, l]$) с погрешностью $O(h^2)$.

Очевидно, что, если вместо уравнения (19.1) рассмотреть уравнение

$$L_1 u := -u''(x) + q(x)u(x) = f(x), \quad x \in (0, 1) \quad (19.21)$$

и аппроксимировать его разностным уравнением

$$L_1^h u^h := -\frac{u_{i+1}^h - 2u_i^h + u_{i-1}^h}{h^2} + q(x_i)u_i^h = f(x_i), \quad i = \overline{1, N-1} \quad (19.22)$$

то задача (19.22), (19.9) будет аппроксимировать задачу (19.21), (19.2) тоже с погрешностью $O(h^2)$.

19.3 Разрешимость и сходимость

Исследуем вопрос о сходимости решения разностной задачи к решению задачи дифференциальной. Для уравнения (19.21) это сделать несколько проще, чем для уравнения (19.1). Поэтому к нему мы и обратимся. Но сначала установим существование и единственность решения задачи (19.22), (19.9).

Теорема 19.1. *Если*

$$q(x) \geq c_1 > 0, \quad 0 < x < 1, \quad (19.23)$$

то решение задачи (19.22), (19.9) существует, единственно, и для него справедлива априорная оценка

$$\max_i |u_i^h| \leq |g_0| + |g_1| + \max_i \frac{|f_i|}{c_1}. \quad (19.24)$$

Доказательство. Задача (19.22), (19.9) представляет собой систему линейных алгебраических уравнений с квадратной матрицей порядка $(N+1)$. Поэтому всегда существует такая правая часть $[g_0, f_1, \dots, f_{N-1}, g_1]$ этой системы (берется первое уравнение из (19.9), затем последовательно все уравнения (19.22) и, наконец, второе уравнение (19.9)), что решение u^h существует. Например, возьмем произвольный набор чисел

$u_0^h, u_1^h, \dots, u_N^h$ и подставим его в левые части (19.22), (19.9). Этим мы определим правые части (19.22), (19.9), при которых решение заведомо существует.

Получим априорную оценку этого решения. Пусть

$$\max_i |u_i^h| = |u_{i_0}^h|.$$

Если $i_0 = 0$ или $i_0 = N$, то в силу (19.9)

$$\max_i |u_i^h| \leq \max\{|g_0|, |g_1|\} \leq |g_0| + |g_1|, \quad (19.25)$$

что согласуется с (19.24). В противном случае максимум модуля достигается во внутреннем узле $x_{i_0} \in \omega$. Запишем уравнение (19.22) в этом узле

$$-\frac{u_{i_0-1}^h - 2u_{i_0}^h + u_{i_0+1}^h}{h^2} + q_{i_0} u_{i_0}^h = f_{i_0}.$$

Если $u_{i_0}^h \geq 0$, то

$$-\left[\underset{0}{\wedge} (u_{i_0-1}^h - u_{i_0}^h) + (u_{i_0+1}^h - \underset{0}{\wedge} u_{i_0}^h) \right] \geq 0$$

и, следовательно,

$$q_{i_0} u_{i_0}^h \leq f_{i_0}.$$

Отсюда с учетом (19.23)

$$0 \leq u_{i_0}^h \leq \frac{f_{i_0}}{q_{i_0}} \leq \frac{1}{c_1} \max_i |f_i|. \quad (19.26)$$

Если же $u_{i_0}^h < 0$, то

$$-\left[(u_{i_0-1}^h - \underset{0}{\vee} u_{i_0}^h) + (u_{i_0+1}^h - \underset{0}{\vee} u_{i_0}^h) \right] \leq 0 \quad (19.27)$$

и, следовательно,

$$q_{i_0} u_{i_0}^h \geq f_{i_0}.$$

Отсюда

$$-q_{i_0} |u_{i_0}^h| \geq f_{i_0}$$

и снова

$$|u_{i_0}^h| \leq -\frac{f_{i_0}}{q_{i_0}} \leq \frac{1}{c_1} \max_i |f_i|. \quad (19.28)$$

Собирая оценки (19.25), (19.26), (19.28), приходим к (19.24). Априорная оценка получена.

Докажем теперь, что решение единственно. Допустим противное, т.е. допустим существование двух решений $u_{(1)}^h$ и $u_{(2)}^h$. Очевидно, что их разность $z = u_{(1)}^h - u_{(2)}^h$ удовлетворяет однородному уравнению (19.22) и однородным граничным условиям (19.9). В силу априорной оценки (19.24)

$$\max_i |z_i| \leq 0.$$

Следовательно, $z_i \equiv 0$, что противоречит предположению. Мы доказали, что однородная система (19.22), (19.9) имеет лишь тривиальное решение. Следовательно, матрица этой системы невырождена, и задача (19.22), (19.9) имеет единственное решение при любых g_0, g_1 и f_i . Теорема доказана.

Теорема 19.2. *Если выполнено условие (19.23), и решение $u(x)$ задачи (19.21), (19.2) принадлежит $C^4[0, l]$, то решение u^h задачи (19.22), (19.9) сходится к решению задачи (19.21), (19.2) со скоростью $O(h^2)$, т.е.*

$$|u(x_i) - u_i^h| = O(h^2).$$

Доказательство. Напишем задачу для погрешности решения $z_i = u_i^h - u(x_i)$. Будем иметь

$$-\frac{z_{i+1} - 2z_i + z_{i-1}}{h^2} + q_i z_i = \Psi_i, \quad z_0 = z_N = 0. \quad (19.29)$$

К задаче (19.29) применим теорему 19.1, в силу которой

$$\max_i |z_i| \leq \frac{1}{c_1} \max_i |\Psi_i|.$$

Но в силу вышедоказанного $\Psi_i = O(h^2)$, что и доказывает теорему.

Замечание 19.3. Более детальный анализ показывает, что

$$\max_i |u_i^h - u(x_i)| \leq \frac{1}{c_1} \max_{x \in [0, l]} |u^{IV}(x)| \frac{h^2}{12}.$$

Теорема 19.3 (О монотонности). *Если выполнено условие*

$$q_i \geq 0, \quad i = \overline{1, N-1}, \quad (19.30)$$

а сеточная функция $U_i, i = \overline{0, N}$ такова, что

$$U_0 \geq 0, \quad U_N \geq 0 \quad (19.31)$$

и

$$L_1^h U_i \geq 0, \quad i = \overline{1, N-1}, \quad (19.32)$$

то

$$U_i \geq 0, \quad i = \overline{1, N-1}. \quad (19.33)$$

Доказательство. Допустим противное, т.е. допустим, что функция U_i может принимать отрицательные значения. Тогда существует такой узел $x_{i_0}, i_0 \in \{1, 2, \dots, N-1\}$, что

$$\min_i U_i = U_{i_0} < 0 \quad (19.34)$$

и в силу (19.27)

$$-(U_{i_0-1} - 2U_{i_0} + U_{i_0+1}) \leq 0.$$

Исследуем обе эти возможности. Если $-(U_{i_0-1} - 2U_{i_0} + U_{i_0+1}) < 0$, то с учетом (19.30) и (19.34)

$$L_1^h U_{i_0} = -U_{i_0-1} - 2U_{i_0} + U_{i_0+1} + q_{i_0} U_{i_0} < 0,$$

и мы пришли к противоречию с (19.32). Если $(U_{i_0-1} - 2U_{i_0} + U_{i_0+1}) = 0$, а $q_{i_0} \neq 0$, мы снова получаем противоречие. Для выхода из этих противоречий мы должны предположить, что $q_{i_0} = 0$ и $(U_{i_0-1} - 2U_{i_0} + U_{i_0+1}) = 0$. Но в силу (19.27), (19.34) это означает, что $U_{i_0-1} = U_{i_0} = U_{i_0+1} < 0$, и в качестве i_0 из (19.34) можно взять также $(i_0 - 1)$ или $(i_0 + 1)$. Делая этот выбор, мы теми же рассуждениями приходим к утверждению, что и $U_{i_0-2} = U_{i_0}$ (или $U_{i_0+2} = U_{i_0}$). И т.д. Поскольку в силу (19.31), (19.34) функция U_i , $i = \overline{0, N}$ не является постоянной, то существует такой узел x_{i_1} , $i_1 \in \{1, 2, \dots, N-1\}$, что $U_{i_1} = U_{i_0}$, а U_{i_1-1} или U_{i_1+1} больше U_{i_1} . В этом узле $-(U_{i_1-1} - 2U_{i_1} + U_{i_1+1}) < 0$, и мы вернулись к уже рассмотренному случаю, который привел нас к противоречию с (19.32). Все противоречия снимаются, если мы откажемся от предположения, что U_i может принимать отрицательные значения. Теорема доказана.

Определение 19.10. Матрица A называется монотонной, если любой вектор x , для которого $Ax \geq 0$, является неотрицательным.

Теорема 19.4 (Принцип сравнения). Пусть u_i^h — решение задачи (19.22), (19.9), а U_i — решение следующей задачи:

$$L_1^h U_i = F_i, \quad i = \overline{1, N-1}, \quad U_0 = G_0, \quad U_N = G_1.$$

Пусть

$$|f_i| \leq F_i, \quad |g_0| \leq G_0, \quad |g_1| \leq G_1. \quad (19.35)$$

Тогда, если выполнено условие (19.30), то

$$|u_i^h| \leq U_i, \quad i = \overline{1, N-1}. \quad (19.36)$$

Доказательство. Легко видеть, что функция $(U_i - u_i^h)$ является решением задачи

$$\begin{aligned} L_1^h (U - u^h)_i &= F_i - f_i, \quad i = \overline{1, N-1}, \quad U_0 - u_0^h = G_0 - g_0, \\ U_N - u_N^h &= G_1 - g_1. \end{aligned}$$

В силу (19.35) и теоремы 19.3 заключаем, что $U_i - u_i^h \geq 0$. Из аналогичных соображений находим, что и $U_i + u_i^h \geq 0$. Тем самым, $-U_i \leq u_i^h \leq U_i$, и теорема доказана.

Замечание 19.4. Функция U_i из (19.36) называется барьером.

Теорема 19.5. Для решения задачи (19.22), (19.9) при выполнении условия (19.30) справедлива априорная оценка

$$\max_i |u_i^h| \leq |g_0| + |g_1| + \frac{l^2}{8} \max_i |f_i|.$$

Доказательство. Введем в рассмотрение функцию

$$U_i = |g_0|(1 - x_i) + x_i|g_1| + cx_i(1 - x_i) \geq 0, \quad (19.37)$$

где $c > 0$ — некоторая постоянная. Очевидно, что $U_0 = |u_0^h|$, $U_N = |u_N^h|$. Легко проверить, что

$$L_1^h U_i = 2c + q_i U_i =: F_i \geq 2c.$$

Пусть $c = 1/2 \max_i |f_i|$. Тогда $|f_i| \leq F_i$, и мы находимся в условиях теоремы 19.4, т.е. $|u_i^h| \leq U_i$. Но

$$\max_i U_i \leq |g_0| + |g_1| + c/4.$$

Теорема доказана.

Упражнение 19.2. Сформулировать и доказать теорему о скорости сходимости разностной задачи (19.22), (19.9).

19.4 Уравнения с переменными коэффициентами

Рассмотрим общее самосопряженное уравнение второго порядка

$$-\frac{d}{dx} \left(p(x) \frac{du}{dx} \right) + q(x)u = f(x), \quad 0 < x < 1. \quad (19.38)$$

и изучим вопрос о его аппроксимации. На первый взгляд кажется вполне естественным раздифференцировать первое слагаемое левой части (19.38)

$$-p(x) \frac{d^2 u}{dx^2} - p'(x) \frac{du}{dx} + q(x)u = f(x) \quad (19.39)$$

и в этом виде заменить $d^2 u/dx^2$ и du/dx соответствующими разностными отношениями. Но так поступать плохо в силу целого ряда причин. В частности, уравнение (19.38) является формально самосопряженным по Лагранжу (симметричным, т.е. если $Lv := -(pv')' + qv$, а $u(x)$ и $v(x)$ обращаются в нуль при $x = 0$ и $x = 1$, то $\int_0^1 vLu dx = \int_0^1 uLv dx$). Сравнить с симметричной матрицей $A = A^T - (Ax, y) = (x, Ay)$). Если же аппроксимировать (19.39), которое эквивалентно (19.38) при гладкой $p(x)$, то аппроксимация, вообще говоря, симметричной не будет. Уравнение (19.38) нужно аппроксимировать сразу в исходном виде.

Построим аппроксимацию (19.38) при помощи интегро - интерполяционного метода (метода баланса, метода конечных объемов). Пусть $x_{i\pm 1/2} = x_i \pm h/2$. Проинтегрируем уравнение (19.38) по отрезку $(x_{i-1/2}, x_{i+1/2})$. Будем иметь

$$\begin{aligned} & -p(x_{i+1/2})u'(x_{i+1/2}) + p(x_{i-1/2})u'(x_{i-1/2}) + \\ & + \int_{x_{i-1/2}}^{x_{i+1/2}} [q(x)u(x) - f(x)] dx = 0. \end{aligned} \quad (19.40)$$

Заменим в (19.40) интеграл квадратурной формулой прямоугольников, а производные — соответствующими разностными отношениями. Именно

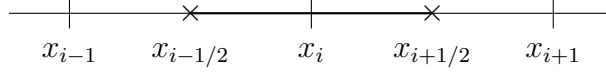


Рис. 1

$$\int_{x_{i-1/2}}^{x_{i+1/2}} [q(x)u - f(x)] dx \approx q_i u_i h - f_i h, \quad (19.41)$$

$$u'_{i+1/2} \approx \frac{u_{i+1} - u_i}{h}, \quad u'_{i-1/2} \approx \frac{u_i - u_{i-1}}{h}.$$

Подставляя (19.41) в (19.40), получим приближенное равенство. Заменяя приближенное равенство на точное, получим уравнение для приближенного решения. После деления на h оно примет вид:

$$-\frac{1}{h} \left[p_{i+1/2} \frac{u_{i+1}^h - u_i^h}{h} - p_{i-1/2} \frac{u_i^h - u_{i-1}^h}{h} \right] + q_i u_i^h = f_i, \quad i = \overline{1, N-1} \quad (19.42)$$

Введем следующие обозначения

$$u_x := u_{x,i} := \frac{u_{i+1} - u_i}{h} \text{ — правое разностное отношение,}$$

$$u_{\bar{x}} := u_{\bar{x},i} := \frac{u_i - u_{i-1}}{h} \text{ — левое разностное отношение.}$$

Очевидно, что $v_{x,i} \equiv v_{\bar{x},i+1}$. Далее,

$$\begin{aligned} \frac{v_{i+1} - 2v_i + v_{i-1}}{h^2} &= \frac{1}{h} \left[\frac{v_{i+1} - v_i}{h} - \frac{v_i - v_{i-1}}{h} \right] = \\ &= \frac{1}{h} (v_{x,i} - v_{\bar{x},i}) = \frac{1}{h} (v_{\bar{x},i+1} - v_{\bar{x},i}) = (v_{\bar{x}})_{x,i} = \\ &= v_{\bar{x}x,i} := v_{\bar{x}x}. \end{aligned}$$

Используя введенные обозначения, уравнение (19.42) можно переписать так:

$$-(p^h u_{\bar{x}}^h)_{x,i} + q_i^h u_i^h = f_i^h, \quad i = \overline{1, N-1}, \quad (19.43)$$

где

$$p^h := p_i^h := p \left(x_i - \frac{h}{2} \right), \quad q^h := q_i^h := q(x_i), \quad f^h := f_i^h := f(x_i). \quad (19.44)$$

19.5 Аппроксимация граничных условий

Применим теперь интегро-интерполяционный метод для построения аппроксимации граничного условия, содержащего производную. Пусть для уравнения (19.38) в точке $x = 0$ (граничной точке) задано граничное условие

$$\alpha \frac{du(0)}{dx} + \beta u(0) = \gamma. \quad (19.45)$$

Граничное условие (19.45) содержит в себе все основные граничные условия для уравнения (19.38): именно, граничные условия первого рода ($\alpha = 0$), второго рода ($\beta = 0$) и третьего рода. Нас будут интересовать граничные условия второго и третьего рода, т.е. условия, содержащие производную. Простейшая аппроксимация условия (19.45) имеет вид

$$\alpha \frac{u_1^h - u_0^h}{h} + \beta u_0^h = \gamma. \quad (19.46)$$

Упражнение 19.3. Доказать, что погрешность аппроксимации граничного условия (19.45) граничным условием (19.46) при $\alpha \neq 0$ есть $O(h)$.

Мы не будем заниматься этой аппроксимацией из-за того, что она имеет большую погрешность. Построим другую аппроксимацию условия (19.45), но прежде его несколько преобразуем. По предположению $\alpha \neq 0$, и на этот коэффициент условие (19.45) можно разделить. Коэффициент $p(x)$ уравнения (19.38) будем предполагать строго положительным

$$p(x) \geq c_0 > 0, \quad (19.47)$$

и домножение (19.45) на $-p(0)$ приведет к эквивалентному уравнению. Будем вместо (19.45) рассматривать граничное условие

$$-p(0) \frac{du(0)}{dx} + \varkappa_0 u(0) = g_0, \quad (19.48)$$

которое при $\alpha = -p(0) \neq 0$, $\beta = \varkappa_0$ и $\gamma = g_0$ совпадает с (19.45). Комбинация $p(0)u'(0)$ в (19.48) хороша уже тем, что величина $-p(x)u'(x)$ имеет смысл потока и фигурирует в самом уравнении (19.38). Знак минус перед производной должен свидетельствовать о том, что производная берется по "внешней нормали": производная $du(0)/dx$ вычислена по направлению внутрь отрезка $[0, 1]$, а производная $-du(0)/dx$ — по направлению, выходящему из отрезка.

Чтобы построить аппроксимацию (19.48), проинтегрируем уравнение (19.38) по отрезку $(0, h/2)$. Будем иметь

$$-p(h/2) \frac{du(h/2)}{dx} + p(0) \frac{du(0)}{dx} + \int_0^{h/2} [q(x)u(x) - f(x)] dx = 0. \quad (19.49)$$

Затем выразим $p(0)du(0)/dx$ из (19.48)

$$p(0)\frac{du(0)}{dx} = \varkappa_0 u(0) - g_0, \quad (19.50)$$

аппроксимируем производную

$$\frac{du(h/2)}{dx} \approx \frac{u_1 - u_0}{h} \quad (19.51)$$

и аппроксимируем интеграл в (19.49) квадратурной формулой "левых прямоугольников"

$$\int_0^{h/2} [q(x)u(x) - f(x)] dx \approx [q(0)u(0) - f(0)]\frac{h}{2}. \quad (19.52)$$

Подставляя теперь (19.50)-(19.52) в (19.49), получим приближенное равенство, которое превратим в точное путем замены точного решения $u(x)$ на приближенное $u^h(x)$. Будем иметь

$$-p_{1/2}\frac{u_1^h - u_0^h}{h} + \left(\varkappa_0 + \frac{h}{2}q_0\right)u_0^h = g_0 + \frac{h}{2}f_0$$

или, принимая обозначения (19.44),

$$-p_1^h u_{\bar{x},1}^h + \left(\varkappa_0 + \frac{h}{2}q_0^h\right)u_0^h = g_0 + \frac{h}{2}f_0^h. \quad (19.53)$$

Соотношение (19.53) представляет собой искомую аппроксимацию.

19.6 Исследование погрешности аппроксимации

Исследуем погрешность аппроксимации разностной схемы (19.43). Исследуем даже более общую схему. Пусть разностная схема имеет вид

$$-\frac{1}{h} [b_i u_{x,i}^h - a_i u_{\bar{x},i}^h] + q_i^h u_i^h = f_i^h. \quad (19.54)$$

Погрешность аппроксимации этой схемы есть

$$\begin{aligned} \Psi_i &= f_i^h + \frac{1}{h} [b_i u_{x,i} - a_i u_{\bar{x},i}] - q_i^h u_i = \\ &= [f_i^h - f(x_i)] - [q_i^h - q(x_i)]u_i + \\ &+ \frac{1}{h} [b_i u_{x,i} - a_i u_{\bar{x},i}] - (pu')'_i. \end{aligned} \quad (19.55)$$

При $u(x) \in C^4[0, 1]$ имеют место следующие разложения

$$\begin{aligned} u_{x,i} &= u'_i + \frac{h}{2}u''_i + \frac{h^2}{6}u'''_i + O(h^3), \\ u_{\bar{x},i} &= u'_i - \frac{h}{2}u''_i + \frac{h^2}{6}u'''_i + O(h^3). \end{aligned}$$

Подставляя эти соотношения в (19.55), будем иметь

$$\begin{aligned}
\Psi_i &= \frac{1}{h} \left[b_i(u'_i + \frac{h}{2}u''_i + \frac{h^2}{6}u'''_i + O(h^3)) - \right. \\
&\quad \left. - a_i(u'_i - \frac{h}{2}u''_i + \frac{h^2}{6}u'''_i + O(h^3)) \right] - \\
&\quad - (p'u' + pu'') - \\
&\quad - [q_i^h - q(x_i)]u_i + [f_i^h - f(x_i)] = \\
&= \left(\frac{b_i - a_i}{h} - p'_i \right) u'_i + \left(\frac{b_i + a_i}{2} - p_i \right) u''_i + \\
&\quad + h \frac{b_i - a_i}{6} u'''_i + O(h^2) - (q_i^h - q_i)u_i + \\
&\quad + (f_i^h - f_i).
\end{aligned}$$

Отсюда находим, что для аппроксимации $O(h^2)$ необходимо и достаточно выполнения условий

$$\begin{aligned}
1^\circ. \quad & \frac{b_i - a_i}{h} - p'_i = O(h^2), \\
2^\circ. \quad & \frac{b_i + a_i}{2} - p_i = O(h^2), \\
3^\circ. \quad & q_i^h - q_i = O(h^2), \\
4^\circ. \quad & f_i^h - f_i = O(h^2).
\end{aligned} \tag{19.56}$$

Для схемы (19.43), (19.44) условия (19.56₃) и (19.56₄) очевидны. Обратимся к (19.56₁) и (19.56₂). Имеем

$$\begin{aligned}
b_i &= p_{i+1/2} = p_i + \frac{h}{2}p'_i + \frac{h^2}{8}p''_i + O(h^3), \\
a_i &= p_{i-1/2} = p_i - \frac{h}{2}p'_i + \frac{h^2}{8}p''_i + O(h^3).
\end{aligned}$$

Отсюда

$$\frac{b_i - a_i}{h} = p'_i + O(h^2), \quad \frac{b_i + a_i}{2} = p_i + O(h^2).$$

Теорема 19.6. *Если решение уравнения (19.38) обладает четвертыми непрерывными производными, то разностная схема (19.43), (19.44) имеет погрешность аппроксимации $O(h^2)$.*

Упражнение 19.4. Доказать, что разностная схема (19.43) при $b_i = a_{i+1}$ и

$$\text{а) } a_i = \frac{p_i + p_{i-1}}{2}, \quad q_i^h = q_i, \quad f_i^h = f_i, \tag{19.57}$$

$$\begin{aligned}
\bar{b}) \quad a_i &= \frac{1}{h} \int_{x_{i-1}}^{x_i} p(x) dx, \quad q_i^h = \frac{1}{h} \int_{x_{i-1}}^{x_{i+1}} q(x)(1 - |x - x_i|) dx, \\
f_i^h &= \frac{1}{h} \int_{x_{i-1}}^{x_{i+1}} f(x)(1 - |x - x_i|) dx
\end{aligned} \tag{19.58}$$

имеет погрешность аппроксимации $O(h^2)$.

Исследуем погрешность аппроксимации ψ_0 граничного условия (19.53). Имеем

$$\begin{aligned}
\psi_0 &:= g_0 + \frac{h}{2} f_0 + p_{1/2} \frac{u_1 - u_0}{h} - (\varkappa_0 + \frac{h}{2} q_0) u_0 = \\
&= g_0 + \frac{h}{2} f_0 + \left(p_0 + \frac{h}{2} p'_0 + O(h^2) \right) \left(u'_0 + \frac{h}{2} u''_0 + O(h^2) \right) - (\varkappa_0 + \frac{h}{2} q_0) u_0 = \\
&= (p_0 u'_0 - \varkappa_0 u_0 + g_0) + \frac{h}{2} (p_0 u''_0 + p'_0 u'_0 - q_0 u_0 + f_0) + O(h^2).
\end{aligned}$$

Первая скобка в этом представлении равна нулю в силу (19.48), а вторая — в силу уравнения (19.38), продолженного по непрерывности с $(0, 1)$ на $[0, 1)$. Тем самым, погрешность аппроксимации граничного условия (19.53) на решении уравнения (19.38) есть $O(h^2)$.

Упражнение 19.5. Интегро-интерполяционным методом построить аппроксимацию граничного условия

$$p(1) \frac{du(1)}{dx} + \varkappa_1 u(1) = g_1 \tag{19.59}$$

и исследовать погрешность полученной аппроксимации.

Теорема 19.7. Пусть выполнены условия

$$p_i^h \geq c_0 > 0, \quad q_i^h \geq c_1 > 0, \quad \varkappa_0 > 0. \tag{19.60}$$

Тогда существует единственное решение задачи (19.43), (19.53), (19.61)

$$u_N^h = g_1, \tag{19.61}$$

и для него справедлива априорная оценка

$$\max_i |u_i^h| \leq \frac{|g_0|}{\varkappa_0} + |g_1| + \max_i \frac{|f_i|}{c_1}. \tag{19.62}$$

Упражнение 19.6. Доказать теорему 19.7.

Теорема 19.8. Если выполнены условия (19.60), и решение задачи (19.38), (19.48), (19.63) $u(x) \in C^4[0, 1]$,

$$u(1) = g_1, \tag{19.63}$$

то решение u^h задачи (19.43), (19.44), (19.53), (19.61) сходится к решению задачи (19.38), (19.48), (19.63) со скоростью $O(h^2)$ равномерно по $x_1 \in \omega$, т.е.

$$\max_i |u(x_i) - u_i^h| = O(h^2).$$

19.7 Некоторые обобщения

Для квазилинейного уравнения

$$-\frac{d}{dx} \left(p(x, u) \frac{du}{dx} \right) + q(x, u) = 0 \quad (19.64)$$

разностную аппроксимацию можно взять в виде

$$\begin{aligned} & -\frac{1}{h} \left[p \left(x_{i+1/2}, \frac{u_{i+1}^h + u_i^h}{2} \right) u_{x,i}^h - p \left(x_{i-1/2}, \frac{u_i^h + u_{i-1}^h}{2} \right) u_{\bar{x},i}^h \right] + \\ & + q(x_i, u_i^h) = 0, \quad i = 1, \dots, N-1. \end{aligned} \quad (19.65)$$

С равным успехом можно поступить и так:

$$\begin{aligned} & -\frac{1}{h} \left[\frac{p(x_{i+1}, u_{i+1}^h) + p(x_i, u_i^h)}{2} u_{x,i}^h - \frac{p(x_i, u_i^h) + p(x_{i-1}, u_{i-1}^h)}{2} u_{\bar{x},i}^h \right] + \\ & + q(x_i, u_i^h) = 0, \quad i = 1, \dots, N-1. \end{aligned} \quad (19.66)$$

Упражнение 19.7. Выяснить порядки погрешности аппроксимации схем (19.65) и (19.66).

Если решения уравнения (19.38) не являются достаточно гладкими, то теорема 19.8 о сходимости со скоростью $O(h^2)$ может не иметь места. В этой ситуации для уменьшения погрешности аппроксимации в окрестности тех точек, где уменьшается гладкость решения, полезно использовать неравномерную сетку. Пусть

$$\widehat{\omega} = \{x_i \mid x_0 = 0 < x_1 < x_2 < \dots < x_{N-1} < x_N = 1\} \quad (19.67)$$

— произвольная неравномерная сетка на $[0, 1]$. Будем обозначать

$$h_i = x_i - x_{i-1}, \quad \bar{h}_i = \frac{h_i + h_{i+1}}{2}.$$

На сетке (19.67) для уравнения (19.38) методом баланса получим следующую аппроксимацию

$$\begin{aligned} & -\frac{1}{\bar{h}_i} \left[p \left(x_i + \frac{h_{i+1}}{2} \right) \frac{u_{i+1}^h - u_i^h}{h_{i+1}} - p \left(x_i - \frac{h_{i+1}}{2} \right) \frac{u_i^h - u_{i-1}^h}{h_i} \right] + \\ & + q(x_i) u_i^h = f(x_i), \quad i = \overline{1, N-1}. \end{aligned} \quad (19.68)$$

Если сетка (19.67) является произвольной, то погрешность аппроксимации (19.68) есть только $O(h)$, где $h = \max h_i$. Однако можно доказать, что погрешность решения соответствующей сеточной задачи и на этой сетке будет величиной $O(h^2)$ при соответствующей гладкости решения уравнения (19.38).

Упражнение 19.8. Исследовать погрешность аппроксимации (19.68).

Выше всюду речь шла о том случае, когда коэффициенты уравнения (19.38) достаточно гладкие. В приложениях часто коэффициенты бывают кусочно-гладкие (например, кусочно-постоянные). В этом случае для аппроксимации уравнения целесообразно использовать сетку, у которой в качестве узлов присутствуют все точки разрыва коэффициентов $p(x)$, $q(x)$ и правой части $f(x)$. Такая сетка будет, как правило, неравномерной и может быть кусочно-равномерной. Указанный выбор сетки позволяет получать точность приближенного решения не ниже, чем в гладком случае.

19.8 Уравнение конвекции-диффузии

Добавим к уравнению еще (19.38) один член — первую производную искомого решения, умноженную на некоторый коэффициент

$$-(pu')' - r(x)u' + q(x)u = f. \quad (19.69)$$

Как аппроксимировать первый и последний члены левой части (19.69), мы знаем. Осталось построить аппроксимацию второго члена. С точки зрения наилучшего порядка аппроксимации следует положить

$$u'_i \approx \frac{u_{i+1} - u_{i-1}}{2h} =: u_{\bar{x}}. \quad (19.70)$$

Тогда аппроксимация уравнения (19.69) примет вид

$$-(p_{i-1/2}u_{\bar{x}})_{x,i} - r_i u_{x,i}^h + q_i u_i^h = f_i, \quad i = 1, \dots, N-1. \quad (19.71)$$

Теорема 19.9. Если $u \in C^4[0, 1]$, то погрешность аппроксимации разностной схемы (19.71) $\Psi = O(h^2)$.

Дополним уравнение (19.69) граничными условиями. Пусть, например,

$$u(0) = g_0, \quad u(1) = g_1. \quad (19.72)$$

Тогда разностные уравнения (19.71) нужно дополнить граничными условиями

$$u_0^h = g_0, \quad u_N^h = g_1. \quad (19.73)$$

Имеет место

Теорема 19.10. Если коэффициенты уравнения (19.69) удовлетворяют условиям (19.23), (19.47), а сетка такова, что

$$\max_i \frac{|r(x_i)|h}{2c_0} \leq 1, \quad (19.74)$$

то задача (19.71), (19.73) имеет единственное решение, и для него справедлива априорная оценка

$$\max_i |u_i^h| \leq |g_0| + |g_1| + \max_i \frac{|f_i|}{c_1}. \quad (19.75)$$

Доказательство. Представим

$$u_{\bar{x}}^h = \frac{u_{i+1}^h - u_{i-1}^h}{2h} = \frac{u_{i+1}^h - u_i^h + u_i^h - u_{i-1}^h}{2h} = \frac{1}{2}u_x + \frac{1}{2}u_{\bar{x}}.$$

Подставим это представление в (19.71)

$$-\frac{1}{h} (p_{i+1/2} u_{x,i}^h - p_{i-1/2} u_{\bar{x},i}^h) - \frac{r_i}{2} (u_{x,i}^h + u_{\bar{x},i}^h) + q_i u_i^h = f_i.$$

Отсюда

$$-\left(\frac{p_{i+1/2}}{h} + \frac{r_i}{2}\right) \frac{u_{i+1}^h - u_i^h}{h} + \left(\frac{p_{i-1/2}}{h} - \frac{r_i}{2}\right) \frac{u_i^h - u_{i-1}^h}{h} + q_i u_i^h = f_i.$$

При выполнении условий (19.74) выражения в скобках неотрицательные. Этого замечания достаточно для того, чтобы завершить доказательство этой теоремы, используя те же самые рассуждения, что и при доказательстве теорем 19.7 и 19.1.

Упражнение 19.9. Завершить доказательство теоремы 19.10.

Упражнение 19.10. Сформулировать и доказать теорему о сходимости разностной задачи (19.71), (19.73).

§ 20

Сингулярно возмущенные уравнения. Негладкие решения

20.1 Осцилляции решения и сингулярно возмущенные уравнения

При исследовании разрешимости и сходимости разностной схемы (19.71) для уравнения конвекции-диффузии (19.69) мы ввели ограничение (19.74) на шаг сетки. Это ограничение в ряде случаев оказывается излишне обременительным, и тогда от аппроксимации (19.70) первой производной приходится отказываться. Обсудим этот вопрос на примере простейшего однородного уравнения с постоянными коэффициентами

$$\frac{d^2 u}{dx^2} + r \frac{du}{dx} = 0, \quad r = \text{const.} \quad (20.1)$$

Наряду с аппроксимацией (19.70) производной u' рассмотрим также ее аппроксимации односторонними разностными отношениями u_x и $u_{\bar{x}}$. Разумеется, порядок погрешности аппроксимации в этих случаях будет хуже. Будем рассматривать одновременно все три из указанных аппроксимаций u' . Для этого в разностное уравнение введем параметр σ

$$u_{\bar{x}\bar{x}}^h + r [\sigma u_x^h + (1 - \sigma) u_{\bar{x}}^h] = 0. \quad (20.2)$$

При $\sigma = 1/2$ имеем $u_x^{\circ} = (u_x + u_{\bar{x}})/2$, при $\sigma = 1 - u_x$, а при $\sigma = 0 - u_{\bar{x}}$. Перепишем (20.2) в поточечном виде

$$\left(\frac{1}{h^2} + \frac{\sigma r}{h} \right) u_{i+1}^h - \left(\frac{2}{h^2} + \frac{(2\sigma - 1)r}{h} \right) u_i^h + \left(\frac{1}{h^2} + \frac{(\sigma - 1)r}{h} \right) u_{i-1}^h = 0.$$

Это есть разностное уравнение с постоянными коэффициентами. Его характеристическое уравнение имеет вид

$$\left(\frac{1}{h} + \sigma r \right) q^2 - \left(\frac{2}{h} + (2\sigma - 1)r \right) q + \left(\frac{1}{h} + (\sigma - 1)r \right) = 0. \quad (20.3)$$

Поскольку сумма коэффициентов уравнения (20.3) равна нулю, то среди его корней есть корень $q_1 = 1$. Второй корень

$$q_2 = q = \frac{1 + (\sigma - 1)rh}{1 + \sigma rh}. \quad (20.4)$$

Проведем качественное сравнение решений дифференциального уравнения (20.1) и разностного уравнения (20.2). Для этого предположим, что

$$r > 0 \quad (20.5)$$

и поставим задачу для (20.1) на положительной полуоси Ox

$$u(0) = 1, \quad u(\infty) = 0. \quad (20.6)$$

Очевидно, что решение задачи (20.1), (20.5), (20.6) имеет вид

$$u(x) = e^{-rx}. \quad (20.7)$$

Функция (20.7) положительна и монотонно убывает при $x \rightarrow \infty$.

Будем искать решение разностного уравнения (20.2), удовлетворяющее условиям

$$u_0^h = 1, \quad \lim_{i \rightarrow \infty} u_i^h = 0. \quad (20.8)$$

Общее решение уравнения (20.2) в силу вышесказанного есть

$$u_i^h = c_1 + c_2 q^i. \quad (20.9)$$

Для того, чтобы это решение на бесконечности было хотя бы ограниченным, нужно потребовать, чтобы (см. (20.4))

$$|q| \leq 1, \quad \text{т.е.} \quad -1 \leq \frac{1 + (\sigma - 1)rh}{1 + \sigma rh} \leq 1. \quad (20.10)$$

Пусть $\sigma \geq 0$. Тогда знаменатель в (20.10) положителен, правая часть неравенства имеет место всегда, и поэтому остается только ограничение

$$-1 - \sigma rh \leq 1 + (\sigma - 1)rh,$$

т.е.

$$2 + (2\sigma - 1)rh \geq 0.$$

Если $\sigma = 1$ или $\sigma = 1/2$, то это условие выполнено со знаком строгого неравенства, и решение задачи (20.2), (20.8) при этих значениях имеет вид

$$u_i^h = q^i, \quad i \in \mathbb{N}. \quad (20.11)$$

Наложим более сильное условие на сеточное решение. Потребуем, чтобы оно было монотонным как и решение (20.7) дифференциальной задачи. Решение (20.11) будет монотонным тогда и только тогда, когда $q \geq 0$, т.е. если

$$1 + (\sigma - 1)rh \geq 0. \quad (20.12)$$

При $\sigma = 1$ это условие выполнено, а при $\sigma = 1/2$ требуется, чтобы

$$rh \leq 2 \quad (20.13)$$

(сравнить с (19.74)).

Итак, если $\sigma = 1$, погрешность аппроксимации уравнения (20.2) есть $O(h)$, но решение (20.11) задачи (20.2), (20.8) монотонно при любых h . Если $\sigma = 1/2$, то погрешность аппроксимации есть $O(h^2)$, но решение (20.11) монотонно только при выполнении условия (20.13). В противном случае решение (20.11) будет колебаться (см. рис. 1), меняя знак при переходе от одного узла к другому.

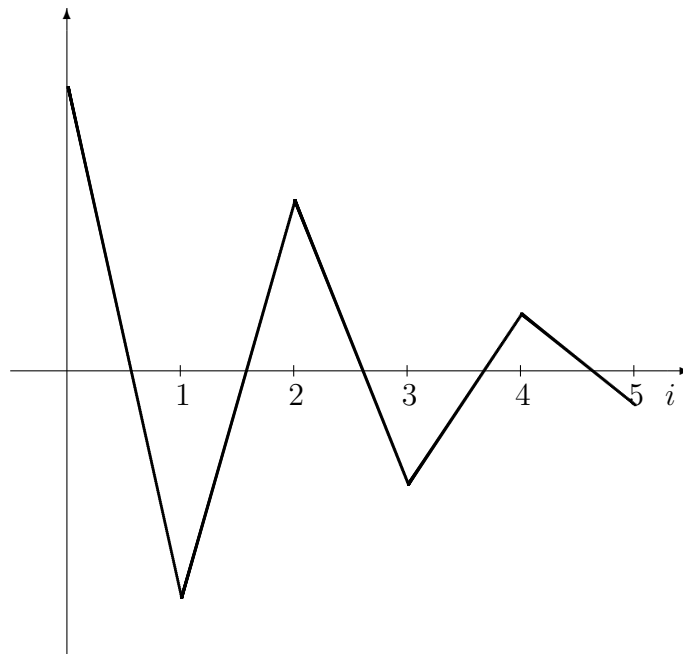


Рис. 1

Именно эти осцилляции решения разностной схемы (20.2) при $\sigma = 1/2$ и не любят прикладники.

Замечание 20.1. Проведенный анализ показал принципиальное различие между схемами (20.2) при $\sigma = 1$ и при $\sigma = 0$, хотя обе эти схемы имеют погрешность $O(h)$ и в этом смысле близки. Причина различия состоит в знаке коэффициента r . Если бы он был отрицательным, то схемы с $\sigma = 1$ и $\sigma = 0$ поменялись бы ролями.

Казалось бы, ограничение (20.13) не является слишком обременительным, чтобы всегда требовать его выполнения. Для обычных задач это так. Но есть важный класс так называемых сингулярно возмущенных уравнений, когда ограничение (20.13) оказывается весьма обременительным. Простейшим примером является уравнение

$$\varepsilon u'' + u' = 0. \quad (20.14)$$

Здесь $\varepsilon \in (0, 1]$ — малый параметр. При $\varepsilon \rightarrow 0$ дифференциальное уравнение второго порядка (20.14) переходит в уравнение первого порядка, для которого одно из двух граничных условий, выделяющих единственное решение уравнения (20.14), становится лишним. Это и является причиной непростого поведения решения соответствующей задачи для уравнения (20.14) при малых ε . Если для уравнения (20.14) поставить граничные условия

$$u(0) = 0, \quad u(1) = 1, \quad (20.15)$$

то решением этой задачи будет функция

$$u(x) = \frac{1 - e^{-x/\varepsilon}}{1 - e^{-1/\varepsilon}} = 1 - \frac{e^{-x/\varepsilon} - e^{-1/\varepsilon}}{1 - e^{-1/\varepsilon}}, \quad (20.16)$$

являющаяся суммой гладкой, медленно меняющейся функции $u_0(x) := 1$ и быстро меняющейся функции $u_1(x) := (e^{-x/\varepsilon} - e^{-1/\varepsilon})/(1 - e^{-1/\varepsilon})$.

Поскольку уравнения (20.1) и (20.14) переходят одно в другое при $r = 1/\varepsilon$, то условие (20.13) применительно к разностной схеме (20.2) для уравнения (20.14) примет вид

$$h \leq 2\varepsilon. \quad (20.17)$$

Но в (20.14) параметр ε может принимать значения 10^{-2} , 10^{-4} или даже 10^{-8} , и ограничение (20.17) становится слишком обременительным. В этой ситуации следует либо ограничиться схемой с $\sigma = 1$, которая не накладывает никаких ограничений на шаг сетки с точки зрения осцилирования решения, и довольствоваться погрешностью аппроксимации $O(h)$, либо пытаться строить другие схемы, которые имеют погрешность $O(h^2)$ и не требуют ограничения типа (20.17).

20.2 Четырехточечная схема

Построим другую аппроксимацию уравнения (20.1). Будем аппроксимировать в (20.1) второе слагаемое при помощи соотношения

$$u'(x_i) \approx \frac{-u_{i+2} + 4u_{i+1} - 3u_i}{2h},$$

погрешность аппроксимации которого есть $O(h^2)$. Используя эту аппроксимацию, вместо (20.2) будем иметь

$$\frac{u_{i+1}^h - 2u_i^h + u_{i-1}^h}{h^2} + r \frac{-u_{i+2}^h + 4u_{i+1}^h - 3u_i^h}{2h} = 0. \quad (20.18)$$

Напишем характеристическое уравнение этого разностного уравнения с постоянными коэффициентами

$$\frac{q^2 - 2q + 1}{h^2} + r \frac{-q^3 + 4q^2 - 3q}{2h} = 0.$$

Обозначим $rh/2 = \xi$ и перепишем характеристическое уравнение в виде

$$-\xi q^3 + (1 + 4\xi)q^2 - (2 + 3\xi)q + 1 = 0.$$

Сумма коэффициентов этого уравнения равна нулю, и следовательно, $q = 1$ есть корень этого уравнения. После деления многочлена из левой части на $(q - 1)$ получим уравнение

$$-\xi q^2 + (1 + 3\xi)q - 1 = 0,$$

корнями которого являются числа

$$q_{2,3} = \frac{1 + 3\xi \pm \sqrt{1 + 6\xi + 9\xi^2 - 4\xi}}{2\xi}.$$

Очевидно, что оба эти корня положительны при любых положительных ξ . Поскольку общее решение уравнения (20.18) имеет вид

$$u_i^h = c_1 + c_2 q_2^i + c_3 q_3^i,$$

то осцилляции этого решения будут отсутствовать на любой сетке, т.е. при любых h .

Аппроксимирующим экспоненту e^{-rh} будет корень

$$\begin{aligned} q_2 &= \frac{1}{2\xi} \left[1 + 3\xi - \sqrt{1 + 2\xi + 9\xi^2} \right] = \\ &= \frac{1}{2\xi} \left\{ 1 + 3\xi - \left[1 + \frac{2\xi + 9\xi^2}{2} - \frac{(2\xi + 9\xi^2)^2}{8} + \frac{8}{16}\xi^3 + O(\xi^4) \right] \right\} = \\ &= 1 - 2\xi + 2\xi^2 + O(\xi^3) = 1 - rh + \frac{r^2 h^2}{2} + O(h^3) = e^{-rh} + O(h^3) \end{aligned}$$

Замечание 20.2. Поскольку уравнение (20.18) нельзя написать для $i = N - 1$, то в этом узле должна быть написана другая аппроксимация уравнения (20.1), например, (20.2) при $i = N - 1$ с любым σ (либо $\sigma = 1/2$, либо $\sigma = 1$).

Замечание 20.3. Мы рассмотрели случай $r > 0$. Если $r < 0$, то, сделав в (20.1) замену независимой переменной $1 - x = t$, придем к уравнению

$$\frac{d^2 u}{dt^2} - r \frac{du}{dt} = u'' + |r|u' = 0.$$

Отсюда следует, что при $r < 0$ в исходных переменных нужно использовать аппроксимацию, зеркальную к той, которая используется при $r > 0$. Именно, вместо разности вперед

$(u_{i+1} - u_i)/h$ — разность назад $(u_i - u_{i-1})/h$, а вместо

$(-u_{i+2} + 4u_{i+1} - 3u_i)/2h$ — аппроксимация

$$u_{\bar{x},i} + \frac{h}{2} u_{\bar{x}x,i} = \frac{u_{i-2} - 4u_{i-1} + 3u_i}{2h}.$$

20.3 О равномерной по ε сходимости

Исследования показывают, что какой бы метод аппроксимации уравнения (20.1) из числа рассмотренных выше мы ни избрали, в любом случае при фиксированном N и $\varepsilon \rightarrow 0$ найдутся такие узлы равномерной сетки, в которых погрешность решения будет $O(1)$. Чтобы отметить этот факт, говорят, что разностная схема не обладает свойством *равномерной по малому параметру сходимости*.

Один из путей обеспечения равномерной по малому параметру сходимости — использование сгущающихся сеток. Одна из простейших сеток, называемая сеткой Шишкина, имеет вид

$$\begin{aligned} \bar{\Omega} = \{x_i \mid x_i = ih, i = \overline{0, N/2}, x_i = x_{N/2} + (i - N/2)H, i = \overline{N/2 + 1, N}, \\ h = \delta/(N/2), \quad H = (1 - \delta)/(N/2), \quad \delta = \min \{c\varepsilon \ln N, 1/2\}\} \end{aligned}$$

или (см. рис. 2)

$$\begin{aligned} x_i = x(t_i), \quad \text{где } t_i = i/N, \quad \text{а} \\ x(t) = \begin{cases} 2\delta t, & 0 \leq t \leq 1/2, \\ 1 - 2(1 - \delta)(1 - t), & 1/2 < t \leq 1 \end{cases} \end{aligned}$$

есть кусочно-линейное непрерывное отображение отрезка $[0, 1]$ на себя.

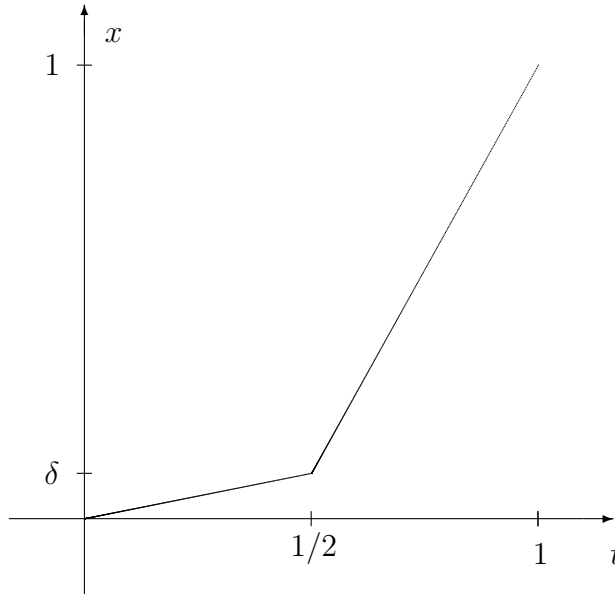


Рис. 2

Эта сетка является кусочно-равномерной с шагом $h \ll H$ на отрезке $[0, \delta]$ и с шагом H на отрезке $[\delta, 1]$.

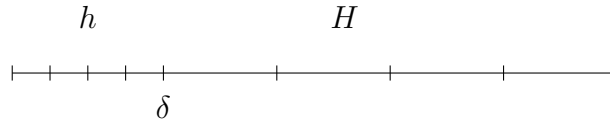


Рис. 3

Равномерная по малому параметру точность разностной схемы определяется погрешностью аппроксимации разностной схемы и величиной параметра ϵ , который должен быть выбран таким, чтобы на длине δ быстро меняющаяся составляющая точного решения успела принять столь малое значение, которое уже не влияет на погрешность приближенного решения.

20.4 Негладкие решения

Рассмотрим следующее дифференциальное уравнение

$$-\frac{1}{x}(xu')' + \frac{\lambda^2}{x^2}u = 0, \quad 0 < x < 1. \quad (20.19)$$

Это уравнение не вкладывается в тот класс уравнений, который мы для себя выделили. Именно, коэффициент $p(x) := x \geq 0$, но не отрезан от нуля постоянной (на рассматриваемом отрезке). Поэтому для уравнения (20.19) в точке $x = 0$ нельзя ставить произвольное граничное условие. В самом деле, будем искать решение уравнения (20.19) в виде

$$u(x) = x^\alpha.$$

Подставляя это выражение в (20.19), находим, что для удовлетворения уравнения требуется выполнение условия

$$\alpha^2 = \lambda^2,$$

т.е. $\alpha = \pm\lambda$. Тем самым, мы нашли два фундаментальных решения уравнения (20.19), и его общее решение есть

$$u(x) = c_1x^\lambda + c_2x^{-\lambda}. \quad (20.20)$$

Без ограничения общности можно считать, что $\lambda > 0$. Если нас интересует ограниченное решение (что естественно с точки зрения приложений), то $c_2 = 0$ и

$$u(x) = c_1x^\lambda.$$

Отсюда находим, что единственным допустимым граничным условием из числа классических является условие

$$u(0) = 0. \quad (20.21)$$

(именно это условие и будет выделять из (20.20) ограниченное решение). При $x = 1$ можно ставить любое граничное условие, например,

$$u(1) = 1. \quad (20.22)$$

Тогда решением задачи (20.19), (20.21), (20.22) будет функция

$$u(x) = x^\lambda. \quad (20.23)$$

Если $0 < \lambda < 1$, то уже первая производная интересующего нас решения не ограничена, не говоря уже о четвертой производной, которая фигурирует в погрешности аппроксимации. О хорошей сходимости численного решения на равномерной сетке говорить трудно. Выход из создавшегося положения можно найти на пути использования специальной сгущающейся к точке $x = 0$ сетки. Как эту сетку построить? Пусть $x = x(t)$ есть отображение отрезка $[0, 1]$ на себя. Для $t \in [0, 1]$ введем равномерную сетку с шагом $h = 1/N$. Тогда

$$x_i = x(t_i)$$

будет задавать узлы неравномерной сетки по x . На этой неравномерной сетке

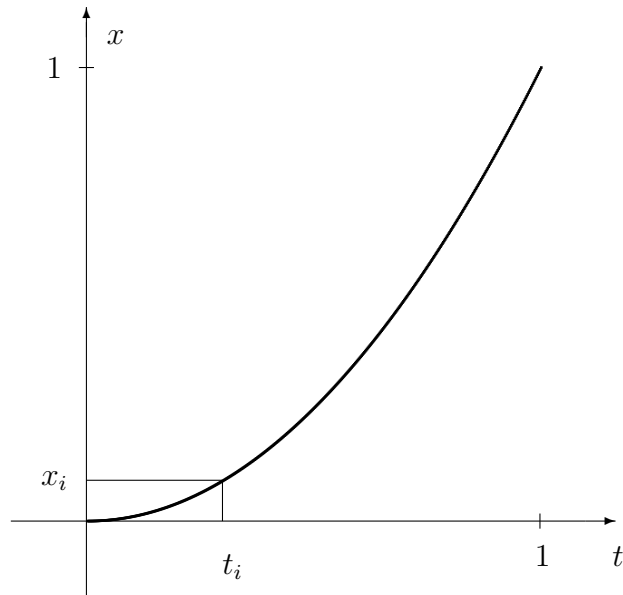


Рис. 4

и аппроксимируем уравнение (20.19). Пусть

$$h_i = x_i - x_{i-1}, \quad \bar{h}_i = (h_i + h_{i+1})/2.$$

Тогда, используя, например, метод баланса (см. § 7), для уравнения (20.19) получим следующую аппроксимацию

$$\frac{1}{x_i} \frac{1}{\bar{h}_i} \left(x_{i+1/2} \frac{u_{i+1}^h - u_i^h}{h_{i+1}} - x_{i-1/2} \frac{u_i^h - u_{i-1}^h}{h_i} \right) - \frac{\lambda^2}{x_i^2} u_i^h = 0. \quad (20.24)$$

Исследуем погрешность аппроксимации этой разностной схемы на неравномерной сетке. Используя формулу Тейлора, будем иметь

$$\begin{aligned} \Psi_i &= \frac{1}{x_i} \frac{1}{\bar{h}_i} [x_{i+1/2} u_{x,i} - x_{i-1/2} u_{\bar{x},i}] - \frac{1}{x_i} (xu')'_i = \\ &= \frac{1}{x_i} \frac{1}{\bar{h}_i} \left[\left(x_i + \frac{h_{i+1}}{2} \right) \left(u'_i + \frac{h_{i+1}}{2} u''_i + \frac{h_{i+1}^2}{6} u'''_i + \frac{h_{i+1}^3}{24} \tilde{u}_i^{IV} \right) - \right. \\ &\quad \left. - \left(x_i - \frac{h_i}{2} \right) \left(u'_i - \frac{h_i}{2} u''_i + \frac{h_i^2}{6} u'''_i - \frac{h_i^3}{24} \tilde{u}_i^{IV} \right) \right] - \frac{1}{x_i} (xu')'_i = \\ &= \frac{1}{x_i} \frac{1}{\bar{h}_i} \left[h_{i+1}^2 \left(\frac{x_i}{6} u'''_i + \frac{1}{4} u''_i \right) + h_{i+1}^3 \left(\frac{x_i}{24} \tilde{u}_i^{IV} + \frac{1}{12} u'''_i \right) + \frac{h_{i+1}^4}{48} \tilde{u}_i^{IV} - \right. \\ &\quad \left. - h_i^2 \left(\frac{x_i}{6} u'''_i + \frac{1}{4} u''_i \right) + h_i^3 \left(\frac{x_i}{24} \tilde{u}_i^{IV} + \frac{1}{12} u'''_i \right) - \frac{h_i^4}{48} \tilde{u}_i^{IV} \right] = \\ &= \frac{h_{i+1}^2 - h_i^2}{\bar{h}_i} \left(\frac{1}{6} u'''_i + \frac{1}{4} \frac{u''_i}{x_i} \right) + \frac{h_{i+1}^3}{\bar{h}_i} \left(\frac{1}{24} \tilde{u}_i^{IV} + \frac{1}{12} \frac{u'''_i}{x_i} \right) + \\ &\quad + \frac{h_i^3}{\bar{h}_i} \left(\frac{1}{24} \tilde{u}_i^{IV} + \frac{1}{12} \frac{u'''_i}{x_i} \right) + \frac{1}{48} \frac{h_{i+1}^4}{\bar{h}_i} \frac{\tilde{u}_i^{IV}}{x_i} - \frac{1}{48} \frac{h_i^4}{\bar{h}_i} \frac{\tilde{u}_i^{IV}}{x_i}. \end{aligned} \quad (20.25)$$

Подставим сюда истинное значение $u(x)$ из (20.23) и оценим вклад в погрешность решения типичной составляющей погрешности аппроксимации

$$\overset{\circ}{\psi}_i = c(x_i) h_i^2 x_i^{\lambda-4}.$$

Эта составляющая представлена в погрешности аппроксимации (20.25) вторым и третьим слагаемыми. Составляющую погрешности решения, отвечающую $\overset{\circ}{\psi}_i$, обозначим через $\overset{\circ}{z}_i$. Для нее имеем уравнение

$$-\frac{1}{x_i} \frac{1}{\bar{h}_i} \left(x_{i+1/2} \frac{\overset{\circ}{z}_{i+1} - \overset{\circ}{z}_i}{h_{i+1}} - x_{i-1/2} \frac{\overset{\circ}{z}_i - \overset{\circ}{z}_{i-1}}{h_i} \right) + \frac{\lambda^2}{x_i^2} \overset{\circ}{z}_i = c(x_i) h_i^2 x_i^{\lambda-4}.$$

Как и при доказательстве теоремы 19.1 для максимума $|\overset{\circ}{z}_i|$ получаем оценку

$$\max_i |\overset{\circ}{z}_i| \leq \max_i \frac{c(x_i) h_i^2 x_i^{\lambda-2}}{\lambda^2}. \quad (20.26)$$

Из этой оценки следует, что, если $\lambda \geq 2$, то никаких проблем нет, ибо в этом случае выражение, стоящее в правой части под знаком \max , имеет равномерную по x_i малость

$O(h_i^2)$, и сетку можно брать равномерной. Если же $\lambda < 2$, то равномерной по x_i малости $O(h_i^2)$ указанного выражения не гарантируется, если сетка не выбрана надлежащим образом. Поскольку $c(x_i)$ из правой части (20.26) меняется мало, выберем сетку при $\lambda < 2$ так, чтобы

$$h_i^2 x_i^{\lambda-2} \approx \text{const.}$$

Так как

$$h_i = x_i - x_{i-1} = N^{-1} x'(t_i^*), \quad (20.27)$$

то

$$h_i^2 x_i^{\lambda-2} = N^{-2} x'^2(t_i^*) x_i^{\lambda-2}.$$

Пусть

$$x'^2 x^{\lambda-2} = c,$$

где c — некоторая постоянная, или

$$x' x^{\lambda/2-1} = \sqrt{c} = c_1.$$

Интегрируя это уравнение, находим, что

$$x^{\lambda/2} = c_1 t + c_2,$$

или

$$x = (c_1 t + c_2)^{2/\lambda}.$$

Так как $x(0) = 0$, а $x(1) = 1$, то $c_2 = 0$, а $c_1 = 1$. Тем самым,

$$x = t^{2/\lambda}, \quad (20.28)$$

и, следовательно,

$$x_i = (i/N)^{2/\lambda}. \quad (20.29)$$

Если узлы сетки будут заданы по закону (20.29), то, в силу (20.27), (20.28) при $\lambda < 2$

$$h_i = 2N^{-1} t_i^{2/\lambda-1} / \lambda,$$

и величины шагов сетки уменьшаются при приближении к границе $x = 0$, т.е. построенная сетка является сгущающейся в окрестности $x = 0$. Если $i \sim N$, то $h_i \sim cN^{-1}$, а если $i = 1$, то

$$h_1 = \frac{2}{\lambda} N^{-2/\lambda} \quad (\lambda < 2).$$

Принимая во внимание сказанное, а также (20.23), (20.28) и (20.26), легко проверить, что вклад последних четырех слагаемых погрешности аппроксимации (20.25) в погрешность решения является величиной $O(N^{-2})$.

Обратимся к первому слагаемому правой части (20.25). Используя, например, формулу Тейлора, находим, что

$$\frac{h_{i+1}^2 - h_i^2}{h_i} = 2(h_{i+1} - h_i) = 2(x_{i+1} - 2x_i + x_{i-1}) = 2N^{-2} x''(t^*).$$

Снова принимая во внимание (20.23), (20.28) и (20.26), заключаем, что вклад и первого слагаемого погрешности аппроксимации (20.25) в погрешность решения оценивается величиной $O(N^{-2})$.

Глава VI

Численные методы для дифференциальных уравнений с частными производными

§ 21

Разностные схемы для уравнения теплопроводности

21.1 Нестационарное уравнение теплопроводности

Нестационарное уравнение теплопроводности является собой простейший пример параболического уравнения — уравнения с частными производными. Возьмем его в виде

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + f(x, t), \quad 0 < x < 1, \quad 0 < t \leq T. \quad (21.1)$$

Чтобы выделить единственное решение уравнения (21.1), нужно задать дополнительные условия. Таковыми могут быть граничные условия, задаваемые при $x = 0$ и $x = 1$, и начальное условие, задаваемое при $t = 0$. Пусть, например, граничные условия имеют вид

$$u(0, t) = u(1, t) = 0, \quad (21.2)$$

а начальное условие —

$$u(x, 0) = \varphi(x). \quad (21.3)$$

Как известно из курса методов математической физики, задача (21.1)-(21.3) поставлена корректно и при надлежащей гладкости $f(x, t)$ и $\varphi(x)$ имеет единственное решение.

Посмотрим на уравнение (21.1) с точки зрения краевых задач для обыкновенных дифференциальных уравнений. Для этого обозначим $\partial u / \partial t = \dot{u}$ и перепишем (21.1) в виде

$$-\frac{\partial^2 u}{\partial x^2} = f(x, t) - \dot{u} \equiv \mathcal{F}(x, t). \quad (21.4)$$

Считая $\mathcal{F}(x, t)$ в (21.4) заданной функцией, а t — параметром, мы можем условно рассматривать (21.4) как обыкновенное дифференциальное уравнение, аппроксимацию которого мы строить умеем. На $[0, 1]$ введем сетку

$$\bar{\omega}^h = \{x = x_i = ih \mid i = 0, \dots, N\}$$

с внутренними узлами

$$\omega^h = \{x_i \in \bar{\omega}^h \mid i = 1, \dots, N-1\}$$

и на этой сетке дифференциальное уравнение (21.4) аппроксимируем разностным уравнением

$$-u_{\bar{x}x,i}^h = \mathcal{F}^h(x_i, t), \quad x_i \in \omega^h. \quad (21.5)$$

Теперь нужно вспомнить (21.4), в силу которого

$$\mathcal{F}(x_i, t) = f(x_i, t) - \dot{u}(x_i, t).$$

Поэтому естественно положить

$$\mathcal{F}^h(x_i, t) = f^h(x_i, t) - \dot{u}_i^h.$$

Тогда (21.5) примет вид

$$-u_{\bar{x}x,i}^h = f_i^h(t) - \dot{u}_i^h, \quad x_i \in \omega^h. \quad (21.6)$$

Это соотношение представляет собой систему $(N-1)$ обыкновенных дифференциальных уравнений первого порядка с $(N+1)$ неизвестными u_i^h , $i = 0, \dots, N$. Воспользуемся граничными условиями (21.2) и положим

$$u_0^h(t) = u_N^h(t) = 0. \quad (21.7)$$

После исключения этих неизвестных из (21.6) будем иметь систему $(N-1)$ уравнений с $(N-1)$ неизвестными.

Перепишем теперь (21.6) по-другому, поставив на первое место производную

$$\dot{u}_i^h = u_{\bar{x}x,i}^h + f_i^h(t), \quad i = 1, \dots, N-1, \quad (21.8)$$

и введем обозначения

$$\begin{aligned} U &= [u_1^h \dots u_{N-1}^h]^T, \\ \Lambda &= \frac{1}{h^2} \begin{bmatrix} -2 & 1 & 0 & \dots & 0 & 0 \\ 1 & -2 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 & -2 \end{bmatrix}, \\ F &= [f_1^h \dots f_{N-1}^h]. \end{aligned} \quad (21.9)$$

Тогда система (21.8) с учетом (21.7) примет вид

$$\frac{dU}{dt} = \Lambda U + F. \quad (21.10)$$

Введем еще одно обозначение

$$\Phi = [\varphi_1 \dots \varphi_{N-1}]^T$$

и положим

$$U(0) = \Phi. \quad (21.11)$$

Соотношения (21.10), (21.11) представляют собой задачу Коши для системы обыкновенных дифференциальных уравнений первого порядка. Для приближенного решения этой задачи можно использовать уже изученные методы. Например, метод Эйлера, который приводит к соотношениям

$$\frac{U^{j+1} - U^j}{\tau} = \Lambda U^j + F^j, \quad U^0 = \Phi, \quad (21.12)$$

или неявный метод Эйлера

$$\frac{U^{j+1} - U^j}{\tau} = \Lambda U^{j+1} + F^{j+1}, \quad U^0 = \Phi, \quad (21.13)$$

а можно и метод трапеций

$$\frac{U^{j+1} - U^j}{\tau} = \frac{1}{2}\Lambda [U^{j+1} + U^j] + \frac{1}{2}(F^{j+1} + F^j), \quad U^0 = \Phi. \quad (21.14)$$

Мы не будем изучать общие методы решения задачи (21.10), (21.11), а ограничимся одношаговыми, как (21.12)-(21.14), которые в теории разностных схем для параболических уравнений принято называть двухслойными.

Изучение (21.12), (21.13) и (21.14) можно проводить одновременно, если записать их единым образом за счет введения параметра σ :

$$\frac{U^{j+1} - U^j}{\tau} = \sigma \Lambda U^{j+1} + (1 - \sigma) \Lambda U^j + \sigma F^{j+1} + (1 - \sigma) F^j. \quad (21.15)$$

Полагая здесь $\sigma = 0, 1$ или $1/2$, получим (21.12), (21.13) или (21.14), соответственно.

Посмотрим теперь на (21.15) с точки зрения аппроксимации не задачи Коши для системы обыкновенных дифференциальных уравнений (21.10), (21.11), а с точки зрения аппроксимации задачи (21.1)-(21.3). В результате двух шагов аппроксимации в области $[0, 1] \times [0, T]$ образована сетка (см. рис. 1),

на которой дифференциальное уравнение (21.1) аппроксимировано системой разностных уравнений

$$\frac{u_i^{hj+1} - u_i^{hj}}{\tau} = \sigma u_{\bar{x},i}^{hj+1} + (1 - \sigma) u_{\bar{x},i}^{hj} + f_i^{hj}, \quad i = 1, \dots, N-1, \quad j = 0, \dots, J-1, \quad (21.16)$$

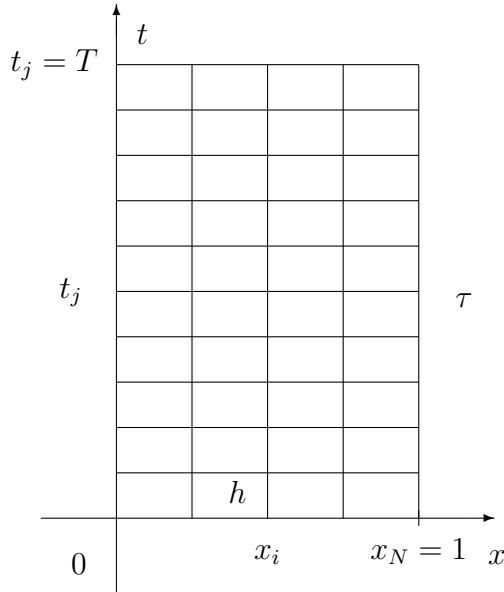


Рис. 1

а граничные (21.2) и начальное (21.3) условия — соотношениями

$$u_0^{hj} = 0, \quad u_N^{hj} = 0, \quad j = 1, \dots, J, \tag{21.17}$$

и

$$u_i^{h0} = \varphi_i, \quad i = 0, \dots, N, \tag{21.18}$$

соответственно. (На связи f_i^{hj} с $f(x, t)$ мы не останавливаемся).

Введем дополнительные обозначения

$$u_i^j = u, \quad u_i^{j+1} = \hat{u}, \quad (\hat{u} - u)/\tau = u_t.$$

В новых обозначениях уравнения (21.16) примут вид

$$u_t^h = \sigma \hat{u}_{\bar{x}\bar{x}}^h + (1 - \sigma)u_{\bar{x}\bar{x}}^h + f^h. \tag{21.19}$$

Погрешностью аппроксимации уравнения (21.1) уравнениями (21.19) будет сеточная функция

$$\Psi = f^h + \sigma \hat{u}_{\bar{x}\bar{x}} + (1 - \sigma)u_{\bar{x}\bar{x}} - u_t,$$

где $u = u(x_i, t_j)$ — значения решения уравнения (21.1) в узлах (x_i, t_j) .

Упражнение 21.1. Доказать, что при надлежащей гладкости (какой?)

$$\Psi = \begin{cases} O(\tau + h^2) & \text{при } \sigma = 0, \sigma = 1, \\ O(\tau^2 + h^2) & \text{при } \sigma = 1/2. \end{cases}$$

Замечание 21.1. Для написания уравнений (21.19) при $\sigma = 0$, $\sigma = 1$ или $\sigma = 1/2$ требуются следующие множества узлов

$$\begin{array}{cccc} j+1 & \bullet & \bullet \bullet \bullet & \bullet \bullet \bullet \\ j & \bullet \bullet \bullet & \bullet & \bullet \bullet \bullet \end{array}$$

соответственно, называемые шаблонами.

21.2 Устойчивость по начальным данным

Исследуем разностную схему (21.16)-(21.18) на предмет ее устойчивости по начальным данным. Для этого будем считать, что правая часть в уравнениях (21.16) равна нулю, т.е.

$$\frac{u_i^{hj+1} - u_i^{hj}}{\tau} = \sigma u_{\bar{x},i}^{hj+1} + (1 - \sigma)u_{\bar{x},i}^{hj}, \quad \begin{array}{l} i = 1, \dots, N - 1, \\ j = 0, \dots, J - 1. \end{array} \quad (21.20)$$

Чтобы исследовать вопрос об устойчивости, найдем решение задачи (21.20), (21.17), (21.18). Решение будем искать методом разделения переменных. Будем искать частные решения уравнений (21.20) в виде

$$u_i^j = X_i T_j.$$

Тогда

$$\frac{T_{j+1} - T_j}{\tau} X_i = (\sigma T_{j+1} + (1 - \sigma)T_j) X_{\bar{x},i}$$

или

$$\frac{(T_{j+1} - T_j)/\tau}{\sigma T_{j+1} + (1 - \sigma)T_j} = \frac{X_{\bar{x},i}}{X_i} = -\lambda^h, \quad (21.21)$$

где λ^h — постоянная. С учетом граничных условий (21.17) для X_i из (21.21) получим задачу

$$X_{\bar{x},i} + \lambda^h X_i = 0, \quad i = 1, 2, \dots, N - 1, \quad X_0 = X_N = 0,$$

или, в развернутом виде,

$$-X_{i-1} + 2X_i - X_{i+1} = h^2 \lambda^h X_i, \quad i = 1, 2, \dots, N - 1, \quad X_0 = X_N = 0. \quad (21.22)$$

Но эта задача совпадает с рассмотренной нами ранее задачей (6.33), если в последней под λ понимать $h^2 \lambda^h$. Поэтому, в силу (6.35)

$$X_i^{(k)} = \sqrt{2} \sin k\pi x_i, \quad i = 1, \dots, N - 1 \quad (21.23)$$

суть собственные векторы задачи (21.22), которые ортогональны в смысле скалярного произведения

$$(u, v) = \sum_{i=1}^{N-1} u_i v_i h$$

и нормированы, т.е. $\|X_i^{(k)}\|^2 = (X_i^{(k)}, X_i^{(k)}) = 1$. В силу (6.42)

$$\lambda_k^h = \frac{4}{h^2} \sin^2 \frac{k\pi h}{2}, \quad k = 1, \dots, N-1 \quad (21.24)$$

— различные собственные значения этой задачи.

Далее, из (21.21) находим, что

$$\frac{T_{j+1}^{(k)} - T_j^{(k)}}{\tau} + \lambda_k^h [\sigma T_{j+1}^{(k)} + (1 - \sigma)T_j^{(k)}] = 0$$

или

$$(1 + \sigma\tau\lambda_k^h)T_{j+1}^{(k)} = (1 - (1 - \sigma)\tau\lambda_k^h)T_j^{(k)}.$$

Отсюда следует, что

$$T_{j+1}^{(k)} = q_k T_j^{(k)},$$

где

$$q_k = \frac{1 - (1 - \sigma)\tau\lambda_k^h}{1 + \sigma\tau\lambda_k^h} \quad (21.25)$$

и поэтому

$$T_j^{(k)} = c_k q_k^j. \quad (21.26)$$

Итак, мы нашли, что функции

$$u_i^{j(k)} = X_i^{(k)} T_j^{(k)}, \quad k = 1, \dots, N-1,$$

где $X_i^{(k)}$ и $T_j^{(k)}$ из (21.23) и (21.26), соответственно, являются частными решениями уравнений (21.20), удовлетворяющими граничным условиям (21.17). Построим линейную комбинацию этих решений

$$u_i^{hj} = \sum_{k=1}^{N-1} c_k X_i^{(k)} q_k^j. \quad (21.27)$$

Полагая здесь $j = 0$, получим

$$u_i^{h0} = \sum_{k=1}^{N-1} c_k X_i^{(k)},$$

а принимая во внимание (21.18), заключаем, что функция (21.27) будет удовлетворять начальным условиям (21.18), если

$$\sum_{k=1}^{N-1} c_k X_i^{(k)} = \varphi_i,$$

т.е. если постоянные c_k суть коэффициенты Фурье функции φ_i при разложении по ортонормированной системе $X_i^{(k)}$

$$c_k = (\varphi, X^{(k)}) = \sum_{i=1}^{N-1} \varphi_i X_i^{(k)} h. \quad (21.28)$$

Итак, сеточная функция (21.27) с коэффициентами c_k из (21.28) удовлетворяет уравнениям (21.20), граничным условиям (21.17) и начальному условию (21.18), а поэтому является решением задачи (21.20), (21.17), (21.18).

Найдем оценку этого решения. Возводя левую и правую части (21.27) в квадрат и суммируя результат по i от 1 до $N-1$, с учетом ортогональности $X_i^{(k)}$, будем иметь

$$\begin{aligned} \|u^{hj}\|_{L_2^h}^2 &= \sum_{i=1}^{N-1} (u_i^{hj})^2 h = \sum_{i=1}^{N-1} h \sum_{k,l=1}^{N-1} c_k c_l X_i^{(k)} X_i^{(l)} q_k^j q_l^j = \\ &= \sum_{k,l=1}^{N-1} c_k c_l q_k^j q_l^j (X_i^{(k)}, X_i^{(l)}) = \sum_{k=1}^{N-1} c_k^2 q_k^{2j} \leq \max_k q_k^{2j} \sum_{k=1}^{N-1} c_k^2 = \max_k q_k^{2j} \|\varphi\|_{L_2^h}^2. \end{aligned}$$

Пусть

$$|q_k| \leq 1. \quad (21.29)$$

Тогда

$$\|u^{hj}\|_{L_2^h} \leq \|\varphi\|_{L_2^h}, \quad (21.30)$$

т.е. L_2^h -норма решения при любом j не превосходит L_2^h -нормы начального условия.

Выясним, когда выполняется условие (21.29). С учетом (21.24), (21.25) при $\sigma \geq 0$ имеем

$$-\left(1 + \frac{4\tau}{h^2} \sigma \sin^2 \frac{k\pi h}{2}\right) \leq 1 - \frac{4\tau}{h^2} (1 - \sigma) \sin^2 \frac{k\pi h}{2}$$

или, после приведения подобных членов,

$$2 - \frac{4\tau}{h^2} (1 - 2\sigma) \sin^2 \frac{k\pi h}{2} \geq 0, \quad k = 1, \dots, N-1.$$

Отсюда вытекает условие

$$(1 - 2\sigma) \leq \frac{h^2}{2\tau} \min_k \frac{1}{\sin^2 \frac{k\pi h}{2}}.$$

Поскольку $\min_k \sin^{-2} \frac{k\pi h}{2} \geq 1$, то (21.29) будет выполнено, если

$$(1 - 2\sigma) \leq \frac{h^2}{2\tau}$$

или, что эквивалентно,

$$\sigma \geq \frac{1}{2} - \frac{h^2}{4\tau}. \quad (21.31)$$

Итак, нами доказана

Теорема 21.1. Если параметр σ схемы (21.20) удовлетворяет условию (21.31), то для решения задачи (21.20), (21.17), (21.18) справедлива априорная оценка

$$\max_j \|u^{hj}\|_{L_2^h} \leq \|u^{h0}\|_{L_2^h}, \quad j = 1, 2, \dots, J. \quad (21.32)$$

Определение 21.1. Говорят, что разностная схема (21.20) устойчива по начальным данным, если для решения задачи (21.20), (21.17), (21.18) справедлива оценка

$$\|u^{hj}\|_{(1)} \leq M \|u^{h0}\|_{(2)},$$

где $\|\cdot\|_{(1)}$ и $\|\cdot\|_{(2)}$ — некоторые нормы, а $M = \text{const} > 0$ не зависит от τ и h .

Следствие 1. Теорема 21.1 утверждает устойчивость по начальным данным схемы (21.20) при выполнении условий (21.31), когда

$$\|\cdot\|_{(2)} = \|\cdot\|_{L_2^h}, \quad \|\cdot\|_{(1)} = \|\cdot\|_{L_\infty(0,T) \times L_2^h(0,1)}.$$

Обсудим условие (21.31). Если $\sigma = 1$, т.е. использован неявный метод Эйлера для системы, то (21.31) выполнено при любых τ и h . То же самое имеет место и при $\sigma = 1/2$ (схема трапеций). Если же $\sigma = 0$, то для выполнения (21.31) нужно, чтобы

$$\tau \leq h^2/2. \quad (21.33)$$

Про первые две схемы (при $\sigma = 1$ и $\sigma = 1/2$) говорят, что они безусловно устойчивы, а третья ($\sigma = 0$) устойчива условно (для устойчивости шага по временной переменной и по пространственной связаны неравенством (21.33)).

Напомним, что все три схемы нуль-устойчивы по терминологии из обыкновенных дифференциальных уравнений, а первые две еще и A -устойчивы.

Отметим, что числа $(-\lambda_k^h)$ являются собственными числами матрицы (21.9)

$$\begin{aligned} \frac{\max_k |-\lambda_k^h|}{\min_k |-\lambda_k^h|} &= \frac{\lambda_{N-1}^h}{\lambda_1^h} = \frac{\sin^2 \frac{(N-1)\pi h}{2}}{\sin^2 \frac{\pi h}{2}} = \\ &= \frac{\cos^2 \frac{\pi h}{2}}{\sin^2 \frac{\pi h}{2}} = \text{ctg}^2 \frac{\pi h}{2} \gg 1 \quad \text{при } h \ll 1, \end{aligned}$$

т.е. система уравнений (21.8) жесткая.

21.3 Устойчивость по правой части

Обратимся теперь к неоднородному уравнению (21.19), а вместо (21.18) поставим однородные начальные условия

$$u_i^{h0} = 0, \quad i = 0, \dots, N. \quad (21.34)$$

Теорема 21.2. Если параметр σ схемы (21.19) удовлетворяет условию (21.31), то для решения задачи (21.19), (21.17), (21.34) справедлива априорная оценка

$$\max_j \|u^{hj}\|_{L_2^h} \leq T \max_j \|f^{hj}\|_{L_2^h}. \quad (21.35)$$

Доказательство. Разложим u_i^{hj} и f_i^{hj} при каждом j по собственным векторам задачи (21.22)

$$u_i^{hj} = \sum_{k=1}^{N-1} T_j^{(k)} X_i^{(k)}, \quad f_i^{hj} = \sum_{k=1}^{N-1} f_j^{(k)} X_i^{(k)}.$$

Подставляя эти разложения в (21.19) и принимая во внимание ортогональность $X_i^{(k)}$, получим

$$\frac{T_{j+1}^{(k)} - T_j^{(k)}}{\tau} + \lambda_k^h [\sigma T_{j+1}^{(k)} + (1 - \sigma) T_j^{(k)}] = f_j^{(k)}.$$

Приводя подобные члены, найдем, что

$$[1 + \sigma\tau\lambda_k^h] T_{j+1}^{(k)} = [1 - (1 - \sigma)\tau\lambda_k^h] T_j^{(k)} + \tau f_j^{(k)},$$

а, разрешая относительно $T_{j+1}^{(k)}$, будем иметь

$$T_{j+1}^{(k)} = q_k T_j^{(k)} + \frac{\tau}{1 + \sigma\tau\lambda_k^h} f_j^{(k)}.$$

В силу (21.29) $|q_k| \leq 1$, а при $\sigma \geq 0$ знаменатель $(1 + \sigma\tau\lambda_k^h) \geq 1$ и поэтому

$$|T_{j+1}^{(k)}| \leq |T_j^{(k)}| |q_k| + \tau |f_j^{(k)}| \leq |T_j^{(k)}| + \tau |f_j^{(k)}|.$$

Далее

$$\|u^{hj+1}\|_{L_2^h} = \sqrt{\sum_{k=1}^{N-1} (T_{j+1}^{(k)})^2} \leq \sqrt{\sum_{k=1}^{N-1} (|T_j^{(k)}| + \tau |f_j^{(k)}|)^2} \leq \|u^{hj}\|_{L_2^h} + \tau \|f^{hj}\|_{L_2^h}.$$

Суммируя это неравенство по j в нужных пределах, приходим к (21.35). Теорема доказана.

Теорема 21.3 (сходимости). Если выполнено условие (21.31), и решение задачи (21.1)-(21.3) $u(x, t) \in C^4[0, 1] \times C^3[0, T]$, то решение u^h задачи (21.16)-(21.18) при соответствующей f_i^{hj} сходится к решению u задачи (21.1)-(21.3) со скоростью $O(h^2 + (\sigma - 1/2)\tau + \tau^2)$.

Доказательство. Пусть $z_i^j = u_i^{hj} - u(x_i, t_j)$ — погрешность решения. Выражая u_i^{hj} через z_i^j и $u(x_i, t_j)$ и подставляя результат в (21.16)-(21.18), для z_i^j получим задачу

$$\begin{aligned} \frac{z_i^{j+1} - z_i^j}{\tau} &= \sigma z_{\bar{x}, i}^{j+1} + (1 - \sigma) z_{\bar{x}, i}^j + \Psi_i^j, \\ z_0^j &= z_N^j = 0, \quad z_i^0 = 0. \end{aligned} \quad (21.36)$$

Для задачи (21.36) справедлива оценка, устанавливаемая теоремой 21.2, т.е.

$$\max_j \|z_i^j\|_{L_2^h} \leq T \max_j \|\Psi_i^j\|_{L_2^h}.$$

Используя теперь результаты упражнения 21.1, приходим к утверждению теоремы.

21.4 Устойчивость в смысле максимума модуля

Теорема 21.4. *Если выполнено условие*

$$\frac{2(1-\sigma)}{h^2}\tau \leq 1, \quad (21.37)$$

то для решения задачи (21.16)-(21.18) справедлива априорная оценка

$$\max_{ij} |u_i^{hj}| \leq \max_i |\varphi_i| + T \max_{ij} |f_i^{hj}|. \quad (21.38)$$

Доказательство. Перепишем уравнение (21.16) в поточечном виде

$$\frac{u_i^{hj+1} - u_i^{hj}}{\tau} = \sigma \frac{u_{i-1}^{hj+1} - 2u_i^{hj+1} + u_{i+1}^{hj+1}}{h^2} + (1-\sigma) \frac{u_{i-1}^{hj} - 2u_i^{hj} + u_{i+1}^{hj}}{h^2} + f_i^{hj}$$

и приведем подобные члены

$$\begin{aligned} & \left(\frac{1}{\tau} + \frac{2\sigma}{h^2} \right) u_i^{hj+1} = \\ & = \frac{\sigma}{h^2} u_{i-1}^{hj+1} + \frac{\sigma}{h^2} u_{i+1}^{hj+1} + \left(\frac{1}{\tau} - \frac{2(1-\sigma)}{h^2} \right) u_i^{hj} + \frac{1-\sigma}{h^2} u_{i-1}^{hj} + \frac{1-\sigma}{h^2} u_{i+1}^{hj} + f_i^{hj}. \end{aligned}$$

Возьмем модули левой и правой частей и оценим правую часть этого соотношения через максимальные значения модулей u_i^{hj} , u_i^{hj+1} и f_i^{hj} . Будем иметь

$$\begin{aligned} & \left(\frac{1}{\tau} + \frac{2\sigma}{h^2} \right) |u_i^{hj+1}| \leq \\ & \leq \frac{2\sigma}{h^2} \max_i |u_i^{hj+1}| + \left(\left| \frac{1}{\tau} - \frac{2(1-\sigma)}{h^2} \right| + \frac{2(1-\sigma)}{h^2} \right) \max_i |u_i^{hj}| + \max_i |f_i^{hj}|. \end{aligned}$$

Беря теперь максимум по i левой части и приводя подобные члены, после домножения на τ получим:

$$\max_i |u_i^{hj+1}| \leq \left(\left| 1 - \frac{2(1-\sigma)\tau}{h^2} \right| + \frac{2(1-\sigma)\tau}{h^2} \right) \max_i |u_i^{hj}| + \tau \max_i |f_i^{hj}|.$$

Изучим коэффициент при $\max_i |u_i^{hj}|$:

$$\left| 1 - \frac{2(1-\sigma)\tau}{h^2} \right| + \frac{2(1-\sigma)}{h^2}\tau = \begin{cases} 1 & \text{при } \frac{2(1-\sigma)}{h^2}\tau \leq 1, \\ \frac{4(1-\sigma)\tau}{h^2} - 1 > 1 & \text{при } \frac{2(1-\sigma)}{h^2}\tau > 1 \end{cases}.$$

В силу условия (21.37) теоремы реализуется первая возможность, и следовательно

$$\max_i |u_i^{hj+1}| \leq \max_i |u_i^{hj}| + \tau \max_i |f_i^{hj}|. \quad (21.39)$$

Пусть

$$\max_j \max_i |u_i^{hj}| = \max_i |u_i^{hj_0}|.$$

Тогда

$$\max_{ij} |u_i^{hj}| = \max_i |u_i^{hj_0}| \leq \max_i |u_i^{hj_0-1}| + \tau \max_i |f_i^{hj_0-1}|.$$

Прибавляя сюда все предыдущие неравенства (21.39) при $j = j_0 - 2, \dots, j = 0$, получим

$$\max_{ij} |u_i^{hj}| \leq \max_i |\varphi_i| + \sum_{j=1}^J \tau \max_i |f_i^{hj-1}| \leq \max_i |\varphi_i| + T \max_{ij} |f_i^{hj}|.$$

Теорема доказана.

Следствие 2. При $\sigma = 1$ оценка (21.38) верна на любой сетке. При $\sigma = 0$ (21.37) совпадает с (21.31), и для справедливости оценки (21.38) должно быть выполнено условие (21.33). Если же $\sigma = 1/2$, то оценка (21.38) верна при

$$\tau/h^2 \leq 1. \quad (21.40)$$

21.5 Сеточное преобразование Фурье

Пусть $\{x_m\}$ — совокупность равноотстоящих узлов на оси Ox . Будем использовать обозначение

$$v(x_m) = v_m, \quad m \in \mathbb{Z}.$$

Будем предполагать, что $v_m \in l_2$, т.е.

$$\sum_{m \in \mathbb{Z}} |v_m|^2 < \infty. \quad (21.41)$$

Определение 21.2. Будем называть 2π -периодическую функцию

$$(Fv_m)(\xi) = \sum_{m \in \mathbb{Z}} v_m e^{-im\xi} = \tilde{v}(\xi) \quad (21.42)$$

сеточным преобразованием Фурье.

Определение 21.3. Обратным сеточным преобразованием Фурье называется сеточная функция

$$(F^{-1}\tilde{v})_m = \frac{1}{2\pi} \int_0^{2\pi} \tilde{v}(\xi) e^{im\xi} d\xi = v_m. \quad (21.43)$$

Замечание 21.2. Соотношение (21.42) на самом деле представляет собой сумму ряда Фурье, коэффициентами которого являются значения рассматриваемой нами сеточной функции v_m . С этой точки зрения (21.43) есть формула для коэффициентов Фурье 2π -периодической функции $\tilde{v}(\xi)$.

Замечание 21.3. Можно было бы называть сеточным преобразованием Фурье функцию

$$Fv_m = \sum_{m \in \mathbb{Z}} hv_m e^{-i(mh)\xi/h} \equiv \tilde{v}(\xi/h), \quad \xi/h = \xi', \quad (21.44)$$

где h — расстояние между соседними узлами $h = x_m - x_{m-1}$. Тогда обратное преобразование приняло бы вид

$$F^{-1}\tilde{v} = \frac{1}{2\pi h} \int_0^{2\pi} \tilde{v}(\xi') e^{i(mh)\xi/h} d\xi = \frac{1}{2\pi} \int_{-\pi/h}^{\pi/h} \tilde{v}(\xi') e^{i(mh)\xi'} d\xi'. \quad (21.45)$$

Устремляя в (21.44) и в (21.45) h к нулю, получим обычные прямое и обратное преобразование Фурье:

$$Fv(x) = \int_{-\infty}^{\infty} v(x) e^{-ix\xi} dx = \tilde{v}(\xi),$$

$$F^{-1}\tilde{v}(\xi) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \tilde{v}(\xi) e^{ix\xi} d\xi = v(x).$$

Для дальнейшего нам потребуется известное из теории рядов Фурье равенство Парсеваля

$$\int_0^{2\pi} |\tilde{v}|^2 d\xi =: \|\tilde{v}\|_{L_2(0,2\pi)}^2 = \frac{1}{2\pi} \|v_m\|_{l_2}^2 := \frac{1}{2\pi} \sum_{m \in \mathbb{Z}} |v_m|^2. \quad (21.46)$$

Пусть T есть оператор сдвига направо, т.е.

$$Tv_m = v_{m+1}.$$

Обратным к нему будет оператор сдвига налево

$$T^{-1}v_m = v_{m-1}.$$

Найдем сеточное преобразование Фурье этих операторов

$$F(Tv_m) = \sum_{m \in \mathbb{Z}} v_{m+1} e^{-im\xi} e^{-i\xi} e^{i\xi} = e^{i\xi} \tilde{v}(\xi).$$

Аналогично

$$F(T^{-1}v_m) = e^{-i\xi} \tilde{v}(\xi).$$

Теперь найдем преобразование Фурье разностных отношений. Имеем

$$Fv_{x,m} = \frac{1}{h} F(T - I)v_m = \frac{e^{i\xi} - 1}{h} \tilde{v}(\xi), \quad (21.47)$$

$$Fv_{\bar{x},m} = \frac{1}{h} (I - T^{-1})v_m = \frac{1 - e^{-i\xi}}{h} \tilde{v}(\xi), \quad (21.48)$$

$$\begin{aligned} Fv_{\bar{x}x,m} &= \frac{1}{h} F(v_{x,m} - v_{\bar{x},m}) = \frac{e^{i\xi} - 1 - 1 + e^{-i\xi}}{h^2} \tilde{v}(\xi) = \\ &= \frac{(e^{i\xi/2} - e^{-i\xi/2})^2}{h^2} \tilde{v}(\xi) = -\frac{4 \sin^2 \xi/2}{h^2} \tilde{v}(\xi). \end{aligned} \quad (21.49)$$

21.6 Устойчивость по начальным данным разностной схемы для уравнения теплопроводности

Рассмотрим разностную схему (21.19) при $f^h \equiv 0$ на сетке, заданной на всей оси Ox , т.е. пусть

$$u_{t,m}^h = \sigma \hat{u}_{\bar{x}x,m}^h + (1 - \sigma) u_{\bar{x}x,m}^h, \quad m \in \mathbb{Z}. \quad (21.50)$$

Теорема 21.5. Если параметр σ схемы (21.50) положителен и удовлетворяет условию

$$\sigma \geq \frac{1}{2} - \frac{h^2}{4\tau} \quad (21.51)$$

то для решения (21.50) имеет место априорная оценка

$$\max_j \|u^{hj}\|_{L_2^h} \leq \|u^{h0}\|_{L_2^h}, \quad j = 1, 2, \dots \quad (21.52)$$

Доказательство. Сделаем в (21.50) сеточное преобразование Фурье

$$\tilde{u}_t + \frac{4 \sin^2 \xi/2}{h^2} (\sigma \hat{u} + (1 - \sigma) \tilde{u}) = 0.$$

Разрешая это обыкновенное разностное уравнение первого порядка относительно \hat{u} , получим

$$\hat{u} = q(\xi) \tilde{u}, \quad (21.53)$$

где

$$q(\xi) = \frac{1 - (1 - \sigma) \frac{4\tau}{h^2} \sin^2 \frac{\xi}{2}}{1 + \sigma \frac{4\tau}{h^2} \sin^2 \frac{\xi}{2}}. \quad (21.54)$$

Из (21.53)

$$\|\hat{u}\|_{L_2(0,2\pi)} = \|q(\xi)\tilde{u}\|_{L_2(0,2\pi)} \leq \max_{0 \leq \xi \leq 2\pi} |q(\xi)| \|\tilde{u}\|_{L_2(0,2\pi)}.$$

Отсюда следует, что L_2 -норма образа Фурье решения не будет возрастать, если

$$|q(\xi)| \leq 1. \quad (21.55)$$

При этом

$$\|\hat{u}\|_{L_2(0,2\pi)} \leq \|\tilde{u}\|_{L_2(0,2\pi)} \leq \dots \leq \|\tilde{u}^0\|_{L_2(0,2\pi)}.$$

Принимая теперь во внимание равенство Парсеваля (21.46), приходим к (21.52).

Покажем теперь, что (21.55) следует из (21.51). Так как $\sigma \geq 0$, то знаменатель в (21.54) положителен, и всегда $q \leq 1$. Осталось проверить условие $q \geq -1$, которое эквивалентно условию

$$2 - (1 - 2\sigma) \frac{4\tau}{h^2} \sin^2 \frac{\xi}{2} \geq 0$$

или

$$1 - 2\sigma \leq \frac{h^2}{2\tau \sin^2 \frac{\xi}{2}}.$$

Но это условие будет выполнено, если

$$1 - 2\sigma \leq \min_{\xi} \frac{h^2}{2\tau \sin^2 \frac{\xi}{2}} = \frac{h^2}{2\tau},$$

что эквивалентно (21.51). Теорема доказана.

Упражнение 21.2. Рассмотреть неоднородное уравнение и установить оценку решения через правую часть.

§ 22

Разностные схемы для уравнения колебаний струны

22.1 Аппроксимация

Рассмотрим другой пример уравнения с частными производными — уравнение колебаний струны

$$\frac{\partial^2 u}{\partial t^2} = \frac{\partial^2 u}{\partial x^2}, \quad 0 < x < 1, \quad 0 < t < T. \quad (22.1)$$

Это — гиперболическое уравнение. Корректной для него является смешанная задача, например,

$$u(0, t) = u(1, t) = 0, \quad u(x, 0) = \bar{u}(x), \quad \frac{\partial u}{\partial t}(x, 0) = \bar{u}'(x). \quad (22.2)$$

Граничные условия при $x = 0$ и $x = 1$ предполагаются однородными граничными условиями первого рода, а в качестве начальных функций взяты некоторые функции $\bar{u}(x)$ и $\bar{u}'(x)$.

Как и при построении аппроксимации уравнения теплопроводности, аппроксимируем сначала производную по пространственной переменной x . В результате получим задачу

$$\begin{aligned} \ddot{u}_i^h(t) &= u_{\bar{x}x, i}^h(t), \quad i = 1, \dots, N-1, \\ u_0^h(t) &= u_N^h(t) = 0, \\ u^h(x_i, 0) &= \bar{u}(x_i), \quad \dot{u}^h(x_i, 0) = \bar{u}'(x_i), \end{aligned} \quad (22.3)$$

которая представляет собой задачу Коши для системы $(N-1)$ дифференциальных уравнений второго порядка. Теперь произведем аппроксимацию по временной переменной: производную $\ddot{u}(t)$ заменим вторым разностным отношением

$$u_{\bar{t}t}(t_j) \equiv [u(t_{j+1}) - 2u(t_j) + u(t_{j-1}))]/\tau^2,$$

а правую часть (22.3) — линейной комбинацией ее значений при $t = t_{j-1}$, $t = t_j$ и $t = t_{j+1}$. В результате будем иметь

$$u_{tt,i}^h = \sigma \hat{u}_{\bar{x}\bar{x},i}^h + (1 - 2\sigma)u_{\bar{x}\bar{x},i}^h + \sigma \check{u}_{\bar{x}\bar{x},i}^h, \quad i = 1, \dots, N - 1, \quad (22.4)$$

где наряду с уже введенным ранее обозначением $\hat{v}_i = v_i(t_{j+1})$ принято обозначение $\check{v}_i = v_i(t_{j-1})$. Правая часть (22.4) представляет собой не общую линейную комбинацию, а линейную комбинацию, симметричную относительно t_{j-1} и t_{j+1} .

К уравнениям (22.4) нужно добавить граничные и начальные условия, которые должны аппроксимировать соответственно условия (22.3)

$$u_0^{hj} = u_N^{hj} = 0, \quad j = 0, \dots, J, \quad (22.5)$$

$$u_i^{h0} = \bar{u}(x_i), \quad i = 1, \dots, N - 1. \quad (22.6)$$

Второе из начальных условий (22.3) содержит производную. Аппроксимируя ее по двум точкам, получим

$$u_{t,i}^{h0} = \bar{u}(x_i), \quad i = 1, \dots, N - 1. \quad (22.7)$$

Теорема 22.1. *Если решение $u(x, t)$ уравнения (22.1) обладает непрерывными четвертыми производными, то погрешность аппроксимации разностной схемы (22.4) есть $O(\tau^2 + h^2)$.*

Доказательство.

$$\begin{aligned} \Psi_i^j &= \sigma \hat{u}_{\bar{x}\bar{x},i}^j + (1 - 2\sigma)u_{\bar{x}\bar{x},i}^j + \sigma \check{u}_{\bar{x}\bar{x},i}^j - u_{tt,i}^j = \sigma \tau^2 u_{\bar{x}\bar{x}t\bar{t},i}^j + u_{\bar{x}\bar{x},i}^j - u_{tt,i}^j = \\ &= \frac{\partial^2 u}{\partial x^2} + O(h^2) - \frac{\partial^2 u}{\partial t^2} + O(\tau^2) = O(\tau^2 + h^2). \end{aligned}$$

Теорема доказана.

Найдем погрешность аппроксимации начального условия (22.7)

$$\psi_i = -u_{t,i}^0 + \bar{u}(x_i) = -\frac{\partial u}{\partial t}(x_i, 0) - \frac{\tau}{2} \frac{\partial^2 u}{\partial t^2}(x_i, 0) + O(\tau^2) + \bar{u}(x_i) = -\frac{\tau}{2} \frac{\partial^2 u}{\partial t^2}(x_i, 0) + O(\tau^2). \quad (22.8)$$

Погрешность аппроксимации начального условия (22.7) есть $O(\tau)$. Построим другую аппроксимацию с погрешностью не хуже $O(\tau^2 + h^2)$. Для этого преобразуем (22.8). В силу (22.1) $\frac{\partial^2 u}{\partial t^2}(x, 0) = \frac{\partial^2 u}{\partial x^2}(x, 0)$, и поэтому

$$\psi_i = -\frac{\tau}{2} \frac{\partial^2 u}{\partial x^2}(x_i, 0) + O(\tau^2).$$

Принимая теперь во внимание (22.2), будем иметь

$$\psi_i = -\frac{\tau}{2} \bar{u}''(x_i) + O(\tau^2).$$

Отсюда и из (22.8) следует, что если вместо $\bar{u}(x_i)$ в (22.7) положить $\bar{u}(x_i) + \frac{\tau}{2}\bar{u}''(x_i)$, т.е. написать условие

$$u_{t,i}^{h0} = \bar{u}_i + \frac{\tau}{2}\bar{u}_i'' \quad (22.9)$$

то погрешность этой аппроксимации будет $O(\tau^2)$. Очевидно также, что если вместо \bar{u}_i'' в (22.9) подставить $\bar{u}_{\bar{x}\bar{x},i}$, т.е. взять

$$u_{t,i}^{h0} = \bar{u}_i + \frac{\tau}{2}\bar{u}_{\bar{x}\bar{x},i} \quad (22.10)$$

то погрешность этой аппроксимации будет $O(\tau^2 + h^2)$.

Аппроксимация (22.10) всем хороша за исключением одного но. Именно, для аппроксимации уравнения (22.1) мы использовали однопараметрическое семейство разностных схем, среди которых имеются как явная ($\sigma = 0$), так и неявные ($\sigma \neq 0$). Аппроксимация же (22.10) всегда явная. Внесем параметр и в начальное условие. Пусть

$$u_{t,i}^{h0} = \bar{u}_i + \frac{\tau}{2}[\gamma u_{\bar{x}\bar{x},i}^{h1} + (1 - \gamma)u_{\bar{x}\bar{x},i}^{h0}].$$

Ясно, что погрешность этой аппроксимации снова не хуже $O(\tau^2 + h^2)$. Наконец, согласуем параметр γ с параметром σ , полагая $\gamma/2 = \sigma$. Аппроксимация второго начального условия (22.2) примет следующий окончательный вид

$$u_{t,i}^{h0} = \tau\sigma u_{\bar{x}\bar{x},i}^{h1} + \tau\left(\frac{1}{2} - \sigma\right)u_{\bar{x}\bar{x},i}^{h0} + \bar{u}_i \quad (22.11)$$

22.2 Устойчивость по начальным данным

Исследуем вопрос об устойчивости схемы (22.4) по начальным данным. Ограничимся изучением задачи Коши, т.е. будем предполагать, что уравнения (22.4) и начальные условия (22.6) и (22.11) заданы для всех $i \in \mathbb{Z}$. Именно, будем рассматривать следующую задачу

$$u_{tt,i}^h = \sigma \hat{u}_{\bar{x}\bar{x},i}^h + (1 - 2\sigma)u_{\bar{x}\bar{x},i}^h + \sigma \check{u}_{\bar{x}\bar{x},i}^h, \quad i \in \mathbb{Z}, \quad (22.12)$$

$$u_i^{h0} = \bar{u}_i, \quad u_{t,i}^{h0} = \tau\sigma u_{\bar{x}\bar{x},i}^{h1} + \tau\left(\frac{1}{2} - \sigma\right)u_{\bar{x}\bar{x},i}^{h0} + \bar{u}_i, \quad i \in \mathbb{Z}. \quad (22.13)$$

Теорема 22.2. *Если параметр σ задачи (22.12), (22.13) неотрицателен и удовлетворяет условию*

$$\sigma \geq \frac{1}{4} - \frac{h^2}{4\tau^2}, \quad (22.14)$$

то для решения этой задачи справедлива априорная оценка

$$\|u^{hj}\|_{L_2^h} \leq \|\bar{u}\|_{L_2^h} + T \|\bar{u}\|_{L_2^h}. \quad (22.15)$$

Доказательство. Сделаем сеточное преобразование Фурье (22.12), (22.13). Для образа Фурье $\tilde{u}^j(\xi)$ решения $u_i^{h,j}$ получим задачу

$$\begin{aligned} \tilde{u}_{\bar{t}\bar{t}} + \frac{4 \sin^2 \frac{\xi}{2}}{h^2} [\sigma \hat{u} + (1 - 2\sigma)\tilde{u} + \sigma \check{u}] &= 0, \\ \tilde{u}^0 = \tilde{\check{u}}, \quad \tilde{u}_t^0 + \frac{4\tau \sin^2 \frac{\xi}{2}}{h^2} \left[\sigma \tilde{u}^1 + \left(\frac{1}{2} - \sigma \right) \tilde{u}^0 \right] &= \tilde{\check{u}}. \end{aligned} \quad (22.16)$$

Умножим теперь уравнение (22.16) на τ^2 и перепишем в поточечном виде

$$\tilde{u}^{j+1} - 2\tilde{u}^j + \tilde{u}^{j-1} + \frac{4\tau^2}{h^2} \sin^2 \frac{\xi}{2} [\sigma \tilde{u}^{j+1} + (1 - 2\sigma)\tilde{u}^j + \sigma \tilde{u}^{j-1}] = 0.$$

Введем обозначение

$$\frac{2\tau}{h} \sin \frac{\xi}{2} = \lambda \quad (22.17)$$

и напишем характеристическое уравнение разностного уравнения

$$(1 + \sigma\lambda^2)q^2 - 2 \left(1 + \left(\sigma - \frac{1}{2} \right) \lambda^2 \right) q + (1 + \sigma\lambda^2) = 0$$

или

$$q^2 - 2 \frac{1 + (\sigma - 1/2)\lambda^2}{1 + \sigma\lambda^2} q + 1 = 0.$$

Отсюда находим корни

$$\begin{aligned} q_{1,2} &= \frac{1 + (\sigma - 1/2)\lambda^2}{1 + \sigma\lambda^2} \pm \frac{\sqrt{[1 + (\sigma - 1/2)\lambda^2 + 1 + \sigma\lambda^2][1 + (\sigma - 1/2)\lambda^2 - 1 - \sigma\lambda^2]}}{1 + \sigma\lambda^2} = \\ &= \frac{1 + (\sigma - 1/2)\lambda^2 \pm \sqrt{[1 + (\sigma - 1/4)\lambda^2](-\lambda^2)}}{1 + \sigma\lambda^2}. \end{aligned} \quad (22.18)$$

Если

$$1 + \left(\sigma - \frac{1}{4} \right) \lambda^2 > 0,$$

то корни q_1 и q_2 будут комплексными и равными по модулю 1. Если же

$$1 + \left(\sigma - \frac{1}{4} \right) \lambda^2 = 0,$$

то

$$q_{1,2} = \frac{1 + (\sigma - 1/2)\lambda^2}{1 + \sigma\lambda^2} = \frac{-1/4\lambda^2}{1/4\lambda^2} = -1$$

и снова $|q_{1,2}| = 1$.

Выясним, когда

$$1 + \left(\sigma - \frac{1}{4} \right) \lambda^2 \geq 0,$$

или, что то же самое,

$$(1/4 - \sigma) \leq \frac{1}{\lambda^2}.$$

Это неравенство будет выполнено при всех $\xi \in (0, 2\pi]$, если (см. (22.17))

$$\frac{1}{4} - \sigma \leq \min_{\xi} \frac{1}{\lambda^2} = \frac{h^2}{4\tau^2}.$$

Но это условие совпадает с условием (22.14) теоремы 22.2, и, следовательно, $|q_{1,2}| = 1$.

Введем обозначение

$$q_{1,2} = e^{\pm i\varphi(\xi)} = \cos \varphi \pm i \sin \varphi,$$

где, согласно (22.18),

$$\cos \varphi = \frac{1 + (\sigma - 1/2)\lambda^2}{1 + \sigma\lambda^2}, \quad \sin \varphi = \frac{|\lambda|\sqrt{1 + (\sigma - 1/4)\lambda^2}}{1 + \sigma\lambda^2}, \quad (22.19)$$

и найдем решение задачи (22.16). Общее решение разностного уравнения (22.16) есть

$$\tilde{u}^j = c_1 \cos j\varphi + c_2 \sin j\varphi.$$

При $j = 0$

$$\tilde{u}^0 = c_1 = \tilde{u},$$

а при $j = 1$

$$\tilde{u}^1 = c_1 \cos \varphi + c_2 \sin \varphi.$$

Из вышесказанного следует, что

$$c_2 = \frac{\tilde{u}^1 - \tilde{u} \cos \varphi}{\sin \varphi}. \quad (22.20)$$

Далее, из второго начального условия (22.16)

$$\tilde{u}^1(1 + \sigma\lambda^2) = \left(1 + \left(\sigma - \frac{1}{2}\right)\lambda^2\right)\tilde{u}^0 + \tau\tilde{u},$$

и, следовательно,

$$\tilde{u}^1 = \frac{1 + (\sigma - 1/2)\lambda^2}{1 + \sigma\lambda^2}\tilde{u}^0 + \frac{\tau\tilde{u}}{1 + \sigma\lambda^2},$$

а с учетом (22.19)

$$\tilde{u}^1 = \cos \varphi \tilde{u} + \frac{\tau\tilde{u}}{1 + \sigma\lambda^2}.$$

Подставляя это значение \tilde{u}^1 в (22.20), найдем, что

$$c_2 = \frac{\tau\tilde{u}}{(1 + \sigma\lambda^2)\sin \varphi}.$$

Окончательно для решения задачи (22.16) получаем представление

$$\tilde{w}^j = \tilde{u} \cos j\varphi + \tau \frac{\tilde{u}}{1 + \sigma\lambda^2} \frac{\sin j\varphi}{\sin \varphi}. \quad (22.21)$$

Чтобы оценить правую часть (22.21), нам потребуется

Лемма 22.1. При $n \in \mathbb{N}$

$$|\sin(n\varphi)/\sin \varphi| \leq n.$$

Доказательство. Имеем

$$\left| \frac{e^{in\varphi} - e^{-in\varphi}}{e^{i\varphi} - e^{-i\varphi}} \right| = \left| \frac{e^{2in\varphi} - 1}{e^{2i\varphi} - 1} \frac{e^{i\varphi}}{e^{in\varphi}} \right| = \left| \frac{e^{2in\varphi} - 1}{e^{2i\varphi} - 1} \right| = |e^{2i(n-1)\varphi} + e^{2i(n-2)\varphi} + \dots + 1| \leq n.$$

Лемма доказана.

Используя лемму 22.1, из (22.21) находим, что

$$|\tilde{w}^j| \leq |\tilde{u}| + \tau j |\tilde{u}| \leq |\tilde{u}| + T |\tilde{u}|.$$

Отсюда

$$\|\tilde{w}^j\|_{L_2(0,2\pi)} \leq \|\tilde{u}\|_{L_2(0,2\pi)} + T \|\tilde{u}\|_{L_2(0,2\pi)},$$

а с учетом равенства Парсеваля

$$\|w^j\|_{L_2^h} \leq \|\bar{u}\|_{L_2^h} + T \|\bar{u}\|_{L_2^h}.$$

Теорема доказана.

Следствие 3. При $\sigma \geq 1/4$ условие (22.14) выполнено, и оценка (22.15) решения имеет место для любых h и τ . При $\sigma = 0$ (явная схема) условие (22.14) выполнено, если $\tau \leq h$, и поэтому оценка (22.15) имеет место только при указанном соотношении между τ и h .