

Оглавление

Глава 1. Численное решение линейных алгебраических систем (СЛАУ).

1. Прямые методы решения СЛАУ. **-1-**
 1. Формулы Крамера.
 2. Метод Гаусса.
 3. Системы с диагональным преобладанием.
 4. Системы с трехдиагональной матрицей. Метод прогонки
2. Обеспечение стабильности СЛАУ. **-10-**
 1. Норм матрицы.
 2. Корректность решения СЛАУ.
 3. Число обусловленности матрицы. Корректность решения СЛАУ.
 4. Оценка числа обусловленности.
 5. Итерационные методы. **-15-**
 1. Построение итерационных последовательностей.
 2. Проблема сходимости итерационного процесса.
 3. Достаточные условия сходимости итерационного процесса.
 4. Метод простой итерации.
 5. Невязочные методы. Метод Зейделя.
 6. Метод верхней релаксации.

Глава 2. Приближение функций.

1. Интерполяция. **-42-**
 1. Классическая постановка задачи интерполяции.
 2. Интерполяция полиномами.
 3. Построение интерполяционного полинома в форме Лагранжа.
 4. Интерполяционный полином в форме Ньютона.
 5. Погрешность интерполяции.
 6. О сходимости интерполяционного процесса.
2. Интерполяция сплайнами. **-55-**
 1. Определение кубического сплайна.
 2. Формулописка системы уравнений для коэффициентов кубического сплайна.
 3. Редукция системы.
 4. Замечание о решении системы.
 5. Сходимость и точность интерполяирования сплайнами.
 6. Метод наименьших квадратов. **-60-**

Глава 1. §1.

- 2 -

$$A_j = \begin{bmatrix} a_{11} & \dots & a_{1,j-1} & f_1 & a_{1,j+1} & \dots & a_{1n} \\ a_{21} & \dots & a_{2,j-1} & f_2 & a_{2,j+1} & \dots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n1} & \dots & a_{n,j-1} & f_n & a_{n,j+1} & \dots & a_{nn} \end{bmatrix} \quad (6)$$

С теоретической точки зрения формулы Крамера (4) дают исчерпывающее решение проблемы. Чтобы найти решение системы (1), нужно подсчитать $n-1$ определитель. Это можно сделать за конечное число арифметических операций. Однако с точки зрения практики важное значение имеет практическое число необходимых операций. Здесь нас и подводят главная трудность. Определитель n -го порядка — это $n!$ слагаемых, каждое из которых является произведением n чисел. Таким образом, для его вычисления нужно выполнить $(n-1)!$ умножений и $n!$ сложений — всего $Q_n = n!$ арифметических операций. Оценим это число. При $n=1$ число $n!$ можно подсчитывать с помощью асимптотической формулы Стирлинга:

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n, \text{ так что } Q_n \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n.$$

При умеренном значении $n=20$ эта формула дает астрономическое число:

$$Q_{20} \approx 5 \cdot 10^{19}.$$

Компьютеру, производительность которого составляет m операций/сек, для вычисления определителя двадцативального порядка понадобится время

$$T_{20} \approx (5 \cdot 10^{19})/m \text{ сек.}$$

В частности, при $m=10^9$ операций/сек получим

$$T_{20} \approx 5 \cdot 10^{-10} \text{ сек.} \approx 170 \text{ лет.}$$

Даже увеличение производительности компьютера на два, три порядка не спасает положения.

Такие результаты получены при $n=20$, в то время, как в современных прикладных задачах приходится решать системы с $n=10^6$ и более уравнений. Из проведенного анализа ясно, что рассчитывать решение СЛАУ по формулам Крамера с вычислением определителей «в лоб» невозможно, т. е. практическая ценность этих формул невелика.

1.2. Метод Гаусса.

Блестящий конструктивный выход из критической ситуации, описанной выше, делает метод Гаусса. Этот метод удобно условно разделить на два этапа. На первом этапе (прямой ход) система (1) приводится к треугольному виду. Затем на втором этапе (обратный ход) осуществляется последовательное отыскание неизвестных x_1, \dots, x_n из этой треугольной системы.

Перед тем как подробно описываем метод Гаусса. Не ограничивая общности, будем считать, что коэффициент a_{11} , который называют ведущим элементом первого шага, отличен от нуля (в случае $a_{11}=0$ поменяем местами уравнения с номерами 1 и i ,

Глава 1. §1.

- 5 -

элементы: $|c_{i,j}| > 1$ и даже $|c_{i,j}| \gg 1$. Тогда при вычислении неизвестных по формулам (12) во время обратного хода умножение найденных с ошибками округления чисел x_i на большие по модулю элементы матрицы C приведет к увеличению этих ошибок. Наоборот, если матрица C оказалась такой, что все ее элементы удовлетворяют условию

$$|c_{i,j}| \leq 1, \quad (17)$$

то роль ошибок округления в процессе вычислений будет нивелироваться.

Опять же, как можно добиться выполнения условия (17). Приступая к первому шагу прямого хода метода Гаусса, рассмотрим элементы a_{ij} первой строки матрицы A и найдем среди них элемент наибольший по модулю. Пусть он имеет номер j_1 . Поменяем в системе (1) первый столбец и столбец с номером j_1 местами, изменяя соответствующим образом нумерацию неизвестных. В результате такой процедуры наибольший по модулю элемент первой строки станет ведущим элементом первого шага a_{11} . Благодаря этому элементы $c_{i,j}$ первой строки матрицы C , которые рассчитываются по формуле (7), будут удовлетворять неравенству (17).

Процедуру выделения наибольшего по модулю элемента в очередной строке и превращение его в ведущий элемент нужно занять повторя во время каждого шага прямого хода метода Гаусса. В этом случае все элементы $c_{i,j}$ треугольной матрицы C (11) будут удовлетворять неравенствам (17), обеспечивая устойчивость метода по отношению к ошибкам округления. Такой способ коррекции называется выбором ведущего элемента по строке.

Появим важность специального выбора ведущего элемента в каждой строке во время прямого хода метода Гаусса на простом примере. Рассмотрим систему трех уравнений с тремя неизвестными:

$$\begin{aligned} 1.2357x_1 + 2.1742x_2 - 5.4834x_3 &= -2.0735 \\ 4.7483x_1 + 6.1365x_2 - 4.7483x_3 &= 4.8755 \\ 6.0969x_1 - 6.2163x_2 + 4.6921x_3 &= 4.8388 \end{aligned} \quad (18)$$

Легко проверить, что ее решение имеет вид

$$x_1 = x_2 = x_3 = 1. \quad (19)$$

Решим систему (18) с помощью метода Гаусса, не обращая внимание на величины элементов матрицы. Все результаты расчетов условимся представлять в виде чисел с плавающей запятой с пятью знаками цифрами. Тогда после прямого хода получим систему треугольного вида:

$$\begin{aligned} x_1 + 1.7595x_2 - 4.4375x_3 &= -1.6780 \\ x_2 + 15324x_3 &= 15324 \\ x_3 &= 0.99992 \end{aligned} \quad (20)$$

Значение $x_3 = 0.99992$ выглядит вполне приемлемым. Однако для двух других неизвестных мы получим следующие значения: $x_2 = 2$, $x_1 = -0.75990$. Причина случившегося заключается в потере точности при вычислении x_2 из-за больших

Глава 3. Численное интегрирование.

1. Формула Ньютона-Лейбница и численное интегрирование. **-65-**
2. Квадратурные формулы прямогольников, трапеций, Симпсона. **-66-**
- 2.1. Сходимость и точность квадратурных формул прямогольников, трапеций и Симпсона.
- 2.2. Апостериорные оценки погрешности при численном интегрировании.
3. Квадратурные формулы Гаусса. **-79-**
- 3.1. Задача построения оптимальных квадратурных формул.
- 3.2. Полиномы Лежандра.
- 3.3. Узлы и весовые коэффициенты квадратурных формул Гаусса.
4. Построение квадратурной формулы с помощью численного интегрирования. **-86-**

Глава 4. Численное интегрирование обыкновенных дифференциальных уравнений.

1. Разностная аппроксимация производных. **-87-**
- 1.1. Сеточные функции.
- 1.2. Разностная аппроксимация первой производной.
- 1.3. Разностная аппроксимация второй производной.
2. Численное решение задачи Коши. **-92-**
- 2.1. Метод Эйтера.
- 2.2. Повышение точности разностного метода.
- 2.3. Метод Рунге-Кутта.
- 2.4. Метод Адамса.
3. Численное решение краевой задачи для линейного дифференциального уравнения второго порядка. **-108-**

Глава 1. ЧИСЛЕННОЕ РЕШЕНИЕ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ СИСТЕМ (СЛАУ)

В этой главе рассматривается одна из самых важных задач линейной алгебры — решение систем линейных алгебраических уравнений, в которых число уравнений равно числу неизвестных:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= f_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= f_2 \\ \dots & \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n &= f_m \end{aligned} \quad (1)$$

или в сокращенной записи:

$$\sum_{j=1}^n a_{ij}x_j = f_i, \quad i = 1, 2, \dots, n.$$

Коэффициенты a_{ij} при неизвестных x_j образуют матрицу системы (1)

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \quad (2)$$

Всюду на протяжении этой главы мы будем считать определитель матрицы отличным от нуля

$$\Delta = \det A \neq 0. \quad (3)$$

В этом случае система (1) называется невырожденной. Решение невырожденной системы всегда существует и является единственным. Обсудим метод фактического построения этого решения.

§1. Прямые методы решения СЛАУ.

Прямыми называются методы, которые позволяют получить точное решение невырожденной системы (1) за конечное число операций.

1.1. Формулы Крамера

Формулы Крамера представляют компоненты x_j решения системы (1) в виде отношения двух определителей:

$$x_j = \frac{\Delta_j}{\Delta}, \quad j = 1, 2, \dots, n, \quad (4)$$

где

$$\Delta_j = \det A_{jj}, \quad j = 1, 2, \dots, n. \quad (5)$$

Здесь матрица A_j получается из матрицы A заменой ее j -го столбца столбцом правых частей системы (1)

Глава 1. §1.

- 2 -

при котором $a_{11} \neq 0$; поскольку система предполагается невырожденной, то такой номер i заведомо найдется).

Разделим все члены первого уравнения на a_{11} и введем в качестве новых коэффициентов c_{ij} , $i=2, \dots, n$ и правой части y_i отношения

$$c_{11} = \frac{a_{12}}{a_{11}}, \quad c_{12} = \frac{a_{13}}{a_{11}}, \quad \dots, \quad c_{1n} = \frac{a_{1n}}{a_{11}}, \quad y_1 = \frac{f_1}{a_{11}}. \quad (7)$$

Вычтем из каждого i -го уравнения системы ($i=2, \dots, n$) первое уравнение умноженное на a_{1i} . Проделав это, мы исключим неизвестное x_1 из всех уравнений, кроме первого. Преобразованная таким образом система (1) примет эквивалентный вид:

$$\begin{aligned} x_1 + c_{12}x_2 + c_{13}x_3 + \dots + c_{1n}x_n &= y_1 \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n &= f_2^{(1)} \\ a_{32}^{(1)}x_2 + a_{33}^{(1)}x_3 + \dots + a_{3n}^{(1)}x_n &= f_3^{(1)} \\ \dots & \\ a_{n2}^{(1)}x_2 + a_{n3}^{(1)}x_3 + \dots + a_{nn}^{(1)}x_n &= f_n^{(1)}. \end{aligned} \quad (8)$$

Значения новых коэффициентов и правых частей системы (8) вычисляются по формуле:

$$a_{ij}^{(1)} = a_{ij} - a_{1j} \frac{a_{11}}{a_{11}}, \quad f_i^{(1)} = f_i - a_{1i} \frac{f_1}{a_{11}}. \quad (9)$$

Естественно выделить из (8) «уокореннную» систему, содержащую $n-1$ уравнение

$$\begin{aligned} a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n &= f_2^{(1)} \\ a_{32}^{(1)}x_2 + a_{33}^{(1)}x_3 + \dots + a_{3n}^{(1)}x_n &= f_3^{(1)} \\ \dots & \\ a_{n2}^{(1)}x_2 + a_{n3}^{(1)}x_3 + \dots + a_{nn}^{(1)}x_n &= f_n^{(1)}. \end{aligned}$$

Продолжая далее процесс исключения, после $(n-1)$ шага получим исходную систему к виду:

$$\begin{aligned} x_1 + c_{12}x_2 + c_{13}x_3 + \dots + c_{1n}x_n &= y_1 \\ x_2 + c_{23}x_3 + \dots + c_{2n}x_n &= y_2 \\ \dots & \\ x_{n-1} + c_{n-1,n}x_n &= y_{n-1} \\ x_n &= y_n \end{aligned} \quad (10)$$

или в матричной форме

$$Cx = y,$$

где матрица C является верхней треугольной матрицей с единицами на главной диагонали

Глава 1. §1.

- 3 -

при котором $a_{11} \neq 0$; поскольку система предполагается невырожденной.

Обратный ход состоит в последовательном определении неизвестных из системы (10) в обратном порядке:

$$\begin{aligned} x_n &= y_n \\ x_{n-1} &= y_{n-1} - c_{n-1,n}x_n \\ x_{n-2} &= y_{n-2} - c_{n-2,n}x_n - c_{n-2,n-1}x_{n-1} \\ \dots & \\ x_1 &= y_1 - c_{12}x_2 - c_{13}x_3 - \dots - c_{1n}x_n \end{aligned} \quad (12)$$

Подсчитав число арифметических операций, которое требуется выполнить при решении СЛАУ по методу Гаусса. Первый шаг прямого хода, согласно формулам (7) и (9), требует n делений и $n(n-1)$ сложений и умножений. Мы считываем деления отдельно, поскольку для компьютера, как и для человека, это более сложная операция.

Перехода последовательно от n к $n-1$ к $n-2$ и т. д. подсчитываем деления

$$Q_1 = n + (n-1) + \dots + 1 = \frac{1}{2}n(n+1), \quad (13)$$

сложений и умножений

$$Q_2 = n(n-1) + (n-1)(n-2) + \dots + 2 \cdot 1 = \frac{1}{3}n(n^2 - 1). \quad (14)$$

Обратный ход, согласно формулам (12), вообще не требует деления, а необходимое число сложений и умножений подсчитывается по формуле

$$Q_3 = 1 + 2 + \dots + (n-1) = \frac{1}{2}n(n-1). \quad (15)$$

Сравнивая (13) и (14) с (15), мы видим, что обратный ход существенно проще прямого. Сумма (14) и (15) даёт общее число сложений и умножений, необходимое для решения СЛАУ по методу Гаусса:

$$Q = Q_1 + Q_2 + Q_3 = \frac{1}{3}n(n-1)\left(n + \frac{5}{2}\right) = \frac{1}{3}n^3 + O(n^2). \quad (16)$$

Оно не идет ни в какое сравнение с числом $n \cdot n!$, которое требуют формулы Крамера при прямом вычислении определителей.

Описанная выше процедура решения системы (1) методом Гаусса может оказаться неустойчивой по отношению к случайным ошибкам, которые неизбежны при компьютерных расчетах в результате округления чисел из-за конечной длины машинного слова. Действительно, предположим, что в процессе приведения системы (1) к треугольному виду (10) у матрицы C (11) образовалась большие по модулю ошибки

Неравенства (24) означают, что в каждой строке матрицы A диагональный элемент выделен: его модуль больше суммы модулей всех остальных элементов той же строки.

Теорема

Система с диагональным преобладанием всегда разрешима и примет единственным образом.

Рассмотрим соответствующую однородную систему:

$$\sum_{j=1}^n a_{ij}x_j = 0, \quad 1 \leq i \leq n \quad (25)$$

Предположим, что она имеет нетривиальное решение \bar{x}_n . Пусть наибольшая по модулю компонента этого решения соответствует индексу $j=k$, т. е.

$$|\bar{x}_k| > |\bar{x}_i| \geq |\bar{x}_j|, \quad 1 \leq j \leq n. \quad (26)$$

Запишем k -ое уравнение системы (25) в виде

$$a_{kk}\bar{x}_k = \sum_{j=1}^n a_{kj}\bar{x}_j \leq |\bar{x}_k| \sum_{j=1}^n |a_{kj}|. \quad (27)$$

Сокращая неравенство (27) на множитель $|\bar{x}_k|$, который, согласно (26), не равен нулю, придем к противоречию с неравенством (24), выражющим диагональное преобладание. Полученное противоречие позволяет последовательно высказать три утверждения:

1. Однородная система (25) с диагональным преобладанием имеет только тривиальное решение.
2. Определитель матрицы A с диагональным преобладанием не равен нулю.
3. Неоднородная система (1) с диагональным преобладанием всегда разрешима и примет единственным образом.

Последнее из них означает, что доказательство теоремы завершено.

1.4. Системы с треугольной матрицей. Метод прогонки.

При решении многих задач приходится иметь дело с системами линейных уравнений виде:

$$A_{ij}x_j + C_{ij}x_{i+1} + B_{ij}x_{i+2} = F_i, \quad i = 1, \dots, n-1, \quad (28)$$

$$x_0 = q_0, \quad x_n = q_n. \quad (29)$$

где коэффициенты A_i, C_i, B_i , правые части F_i ($i = 1, \dots, n-1$) известны вместе с числами q_i и q_n . Дополнительные соотношения (29) часто называют краевыми условиями для системы (28). Во многих случаях они могут иметь более сложный вид. Например:

$$x_0 = p_0 x_1 + q_0, \quad x_n = p_n x_{n-1} + q_n,$$

где p_0, q_0, p_n, q_n – заданные числа. Однако, чтобы не усложнять изложение, мы ограничимся простейшей формой дополнительных условий (29).

Пользуясь тем, что значения x_0 и x_n заданы, перенесем систему (28) в виде:

$$\begin{aligned} C_1 x_1 + B_1 x_2 &= F_1 - A_1 q_0 \\ A_2 x_1 + C_2 x_2 + B_2 x_3 &= F_2 - B_1 q_0 \\ \vdots & \\ A_{n-1} x_{n-2} + C_{n-1} x_{n-1} &= F_{n-1} - B_{n-2} q_n \end{aligned} \quad (30)$$

Матрица этой системы имеет трехдиагональную структуру:

$$\begin{bmatrix} C_1 & B_1 & 0 & \cdots & 0 & 0 \\ A_2 & C_2 & B_2 & 0 & \cdots & 0 \\ 0 & A_3 & C_3 & B_3 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & 0 & \cdots & A_{n-1} & C_{n-1} \end{bmatrix} \quad (31)$$

Это существенно упрощает решение системы (28) благодаря специальному методу, полученному названию метода прогонки.

Метод основан на предположении, что искомые неизвестные x_i и x_{i+1} связаны рекуррентным соотношением

$$x_i = \alpha_i x_{i+1} + \beta_{i+1}, \quad 0 \leq i \leq n-1. \quad (32)$$

Здесь величины α_{i+1} , β_{i+1} , получившие название прогоночных коэффициентов, подлежат определению, исходя из условий задачи (28), (29). Фактически такая процедура означает замену прямого определения неизвестных x_i задачей определения прогоночных коэффициентов с последующим расчетом по ним величин x_i .

Для реализации описанной программы выразим с помощью соотношения (32) x_{i+1} через x_{i+1} :

$$x_i = \alpha_i x_{i+1} + \beta_{i+1} = \alpha_i \alpha_{i+1} x_{i+2} + \alpha_i \beta_{i+1} + \beta_i, \quad (33)$$

и подставим x_{i+1} и x_i , выраженные через x_{i+1} , в исходные уравнения (28). В результате получим:

$$(A_i \alpha_i \alpha_{i+1} + C_i \alpha_{i+1} + B_i \beta_{i+1}) x_{i+1} + A_i \alpha_i \beta_{i+1} + A_i \beta_i + C_i \beta_{i+1} = F_i = 0, \quad i = 1, 2, \dots, n-1.$$

Последние соотношения будут заведомо выполняться и притом независимо от решения, если потребовать, чтобы при $i = 1, 2, \dots, n-1$ имели место равенства:

$$A_i \alpha_i \alpha_{i+1} + C_i \alpha_{i+1} + B_i \beta_{i+1} = 0,$$

$$A_i \alpha_i \beta_{i+1} + A_i \beta_i + C_i \beta_{i+1} = F_i = 0.$$

Отсюда следуют рекуррентные соотношения для прогоночных коэффициентов:

Отсюда следует, что в конечномерном пространстве норма матрицы ограничена, причем на единичной сфере всегда найдется такой вектор \mathbf{x}_0 , что

$$\|\mathbf{A}\| = \|\mathbf{A}\mathbf{x}_0\|.$$

Наконец, из определения нормы (42) следует, что

$$\|\mathbf{A}\| \leq \|\mathbf{A}\| \|\mathbf{x}_0\|. \quad (44)$$

Это простое неравенство лежит в основе всех дальнейших оценок.

2.2. Корректность решения СЛАУ.

Следуя Адамару, будем называть математическую задачу корректной, если выполняются три условия:

1. Решение задачи существует.

2. Решение задачи единственное.

3. Решение задачи непрерывно зависит от входных данных.

Обсудим с точки зрения этого определения задачу решения СЛАУ с неравным нулю определителем

$$A\mathbf{x} = \mathbf{f}, \quad (45)$$

считая матрицу A фиксированной и рассматривая в качестве входных данных вектор правых частей системы $\mathbf{f} = [f_1, f_2, \dots, f_n]^T \in E_n$.

Условие $\Delta \neq 0$ гарантирует существование у матрицы A обратной матрицы A^{-1} , через которую решение системы (45) можно записать в виде

$$\mathbf{x} = A^{-1} \mathbf{f}. \quad (46)$$

Пусть первая правая часть подверглась возмущению $\delta\mathbf{f}$ и стала равной $\tilde{\mathbf{f}} = \mathbf{f} + \delta\mathbf{f}$. Тогда, согласно (46), решение $\tilde{\mathbf{x}}$ возмущенной системы

$$\tilde{\mathbf{x}} = A^{-1} \tilde{\mathbf{f}} = A^{-1} \mathbf{f} + A^{-1} \delta\mathbf{f} = \mathbf{x} + \delta\mathbf{x}, \quad (47)$$

где

$$\delta\mathbf{x} = A^{-1} \delta\mathbf{f}. \quad (48)$$

Отсюда получаем

$$\|\delta\mathbf{x}\| \leq \|A^{-1}\| \|\delta\mathbf{f}\|. \quad (50)$$

Неравенство (50) доказывает непрерывную зависимость возмущения решения $\delta\mathbf{x}$ от возмущения правой части $\delta\mathbf{f}$:

$$\|\delta\mathbf{x}\| \leq \|\delta\mathbf{f}\| \rightarrow 0. \quad (51)$$

Это означает, что решение СЛАУ с неравным нулю определителем Δ -корректная математическая задача: для нее выполняются все три требования корректности Адамара.

2.3. Число обусловленности матрицы.

Исходное уравнение (45) позволяет написать неравенство:

$$\|\mathbf{f}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|. \quad (52)$$

Перемножая его с неравенством того же знака (50), получим:

Она согласуется с результатом (58), который мы получили, непосредственно решая системы (56) и (57).

В процессе решения задачи мы убедились в том, что подсчет числа обусловленности является сложной задачей, особенно с учетом того, что нужно вычислять норму не только прямой, но и обратной матрицы. Поэтому желательно получить какие-нибудь конструктивные оценки этой важнейшей характеристики системы.

2.4. Оценка числа обусловленности.

Исходное уравнение (45) позволяет написать неравенство:

$$\|\mathbf{f}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|. \quad (53)$$

Приемночая е с неравенством того же знака (50), получим:

$$\|\mathbf{f}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \|\delta\mathbf{f}\|. \quad (54)$$

$$\|\mathbf{f}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \|\delta\mathbf{f}\|.$$

Аналогичным образом для собственного вектора \mathbf{z} , связанного с λ_{\max} , имеем

$$\mathbf{A} \mathbf{z} = \lambda_{\max} \mathbf{z},$$

или

$$\mathbf{A}^{-1} \mathbf{z} = \frac{1}{\lambda_{\max}} \mathbf{z}.$$

Отсюда следует оценка

$$\frac{1}{\lambda_{\max}} \leq \|\mathbf{A}^{-1}\|.$$

Перемножая два последних неравенства, приходим к утверждению (60).

Если матрица симметричная $A = A^T$, то все ее характеристические значения вещественны, причем

$$\|\mathbf{A}\| = |\lambda_{\max}| \quad \text{и} \quad \|\mathbf{A}^{-1}\| = \frac{1}{|\lambda_{\min}|}.$$

поскольку для таких матриц

$$\|\mathbf{A}\| = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}. \quad (61)$$

Из полученной оценки для M_A следуют два важных вывода:

$$1) M_A \geq 1;$$

$$2) M_A \leq \|\mathbf{A}\| \|\mathbf{x}\|.$$

$$\alpha_{ii} = \frac{-B_i}{A_i C_i + C_i}, \quad \beta_{ii} = \frac{F_i - A_i \beta_{i-1}}{A_i C_i + C_i}, \quad i = 1, 2, \dots, n-1. \quad (33)$$

Левое граничное условие $x_0 = q_0$ и соотношение $x_0 = \alpha_1 x_1 + \beta_1$ непротиворечивы, если положить

$$\alpha_1 = 0, \quad \beta_1 = q_0. \quad (34)$$

Остальные значения коэффициентов прогонки $\alpha_2, \dots, \alpha_n$ и β_2, \dots, β_n находим из (33), чем и завершается этап вычисления прогоночных коэффициентов.

Далее, согласно правому граничному условию

$$x_n = q_n. \quad (35)$$

Отсюда можно найти остальные неизвестные x_{n-1}, \dots, x_1 в процессе обратной прогонки с помощью рекуррентной формулы (32).

Число операций, которое требуется для решения системы общего вида (1) методом Гаусса, растет при увеличении n пропорционально n^2 . Метод прогонки сводится к двум циклам: сначала по формулам (33) рассчитываются прогоночные коэффициенты, затем с их помощью по рекуррентным формулам (32) находятся компоненты решения системы \mathbf{x} . Это означает, что с увеличением размеров системы число арифметических операций будет расти пропорционально n , а не n^2 . Таким образом, метод прогонки в пределах сферы своего возможного применения является существенно более экономичным. К этому следует добавить особую простоту его программной реализации на компьютере.

Во многих прикладных задачах, которые приводят к СЛАУ с трехдиагональной матрицей, ее коэффициенты удовлетворяют неравенствам:

$$|C_i| > |A_i| + |B_i|, \quad (36)$$

которые выражают свойство диагонального преобладания. В частности, мы встретим такие системы в третьей и пятой главе.

Согласно теореме предыдущего раздела решение таких систем всегда существует и является единственным. Для них также справедливо утверждение, которое имеет важное значение для фактического расчета решения с помощью метода прогонки.

Лемма

Если для системы с трехдиагональной матрицей выполняется условие диагонального преобладания (36), то прогоночные коэффициенты удовлетворяют неравенству:

$$|\alpha_i| \leq 1. \quad (37)$$

Доказательство проведем по индукции. Согласно (34) $\alpha_1 = 0$, т. е. при $i = 1$ утверждение леммы верно. Допустим теперь, что оно верно для α_i , и рассмотрим α_{i+1} :

$$|\alpha_{i+1}| = \left| \frac{B_i}{C_i + A_i \alpha_i} \right| \leq \frac{|B_i|}{|C_i| - |A_i|} \leq 1. \quad (38)$$

Итак, индукция при $i + 1$ обоснована, что и завершает доказательство леммы.

Неравенство (37) для прогоночных коэффициентов α_i делает прогонку устойчивой. Действительно, предположим, что компоненты решения \mathbf{x}_i в результате процедуры округления рассчитана с некоторой ошибкой. Тогда при вычислении

$$\|\delta\mathbf{x}\| \leq \|\delta\mathbf{f}\| \|\mathbf{A}\| \|\mathbf{x}\|. \quad (53)$$

Пусть $\mathbf{f} \neq 0$, тогда, согласно (46), $\mathbf{x} \neq 0$ и неравенство (53) можно переписать в виде:

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq M_A \frac{\|\delta\mathbf{f}\|}{\|\mathbf{f}\|}, \quad (54)$$

где

$$M_A = \|\mathbf{A}\| \|\mathbf{x}\|. \quad (55)$$

Число M_A называется числом обусловленности матрицы A . Оно позволяет оценить относительную погрешность решения через относительную погрешность возмущения правой части. Поскольку исходная система (45) линейная, оценка относительной погрешности является более естественной, чем оценка абсолютной погрешности. Чем больше M_A , тем реальнее реагирует решение на возмущение правой части. Поэтому матрицы с большим числом обусловленности и соответствующие им СЛАУ называют плохо обусловленными. Для оценки роли, которую играет число обусловленности при решении линейных алгебраических систем, разберем задачу.

Задача 1

Рассмотреть систему двух уравнений

$$\begin{aligned} x_1 + 0 \cdot x_2 &= 1 - d \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \\ x_1 + 0.01 \cdot x_2 &= 1' - d \begin{bmatrix} 1 & 0.01 \\ 0 & 1 \end{bmatrix}, \end{aligned} \quad (56)$$

и соответствующую ей возмущенную систему

$$\begin{aligned} x_1 + 0 \cdot x_2 &= 1, \\ x_1 + 0.01 \cdot x_2 &= 1.01, \end{aligned} \quad (57)$$

Выписать решения этих систем, подсчитать погрешность возмущения правой части и соответствующую ей погрешность возмущения решения. Найти число обусловленности матрицы A , составить с его помощью теоретическую оценку погрешности (54) и сравнить результатом, полученным непосредственно по известным решениям систем.

В данном случае определять матрицы A отличнее от нуля.

$$\Delta = \det A \neq 0,$$

т. е. обе системы невырожденные. Система (57) отличается от системы (56) возмущением правой части

$$\mathbf{f} = \{1\}, \quad \|\mathbf{f}\| = \sqrt{2}, \quad \tilde{\mathbf{f}} = \{1, 1.01\}, \quad \|\tilde{\mathbf{f}}\| = 0.001, \quad \|\delta\mathbf{f}\| = 0.001.$$

Решения систем (56) и (57) имеют вид:

$$\mathbf{x} = \{1, 0\}, \quad \|\mathbf{x}\| = 1, \quad \tilde{\mathbf{x}} = \{1, 1\}, \quad \|\tilde{\mathbf{x}}\| = 0.1, \quad \|\delta\mathbf{x}\| = 1.$$

При этом

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \frac{0.1}{\sqrt{2}}, \quad \frac{\|\delta\mathbf{x}\|}{\|\tilde{\mathbf{x}}\|} = 1. \quad (58)$$

Мы видим, что небольшое относительное возмущение правой части привело к сильному возмущению решения: относительная погрешность решения равна единице. Этот результат означает, что исходная система плохо обусловлена. Чтобы убедиться в

приемлемости этой оценки, воспользуемся теоретической оценкой (54).

Число обусловленности

$$M_A = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}.$$

Возьмем для \mathbf{A} единичную матрицу $A = I_2$. Тогда

$$\|\mathbf{A}\| = \|\mathbf{f}\| = \sqrt{2}, \quad \|\mathbf{x}\| = \|\tilde{\mathbf{x}}\| = 1.$$

Следовательно

$$M_A = \sqrt{2}.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

Следовательно

$$\|\delta\mathbf{x}\| = \sqrt{2} \|\mathbf{x}\|.$$

Следовательно

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} = \sqrt{2}.$$

$$\lim_{k \rightarrow \infty} z_k = 0, \quad \lim_{k \rightarrow \infty} z_k^i = 0, \quad 1 \leq i \leq n. \quad (70)$$

Вторая характеристика – невязка:

$$\psi_k = Ax_k - f. \quad (71)$$

Она показывает, насколько хорошо или, наоборот, плохо член итерационной последовательности x_k удовлетворяет исходной системе.

Установим связь между z_k и ψ_k :

$$\psi_k = Ax_k - f = A(z_k + x) - f = Az_k. \quad (72)$$

Можно также написать обратное соотношение:

$$z_k = A^{-1}\psi_k. \quad (73)$$

Из формул (72) и (73) вытекают оценки:

$$\|\psi_k\| \leq \|A\| \|z_k\|, \quad \|z_k\| \leq \|A^{-1}\| \|\psi_k\|. \quad (74)$$

Они показывают, что погрешность z_k стремится к нулю тогда и только тогда, когда стремится к нулю невязка ψ_k . Этот результат позволяет судить о сходимости или расходности итерационного процесса по поведению невязки, которая доступна прямому вычислению и благодаря этому может контролироваться.

При исследовании сходимости итерационных методов большую роль играют свойства матрицы A и B_{n+1} , в первую очередь как самосопряженность и знакоопределенность. Напомним, что в вещественном евклидовом пространстве E_n для каждого линейного преобразования существует единственное сопряженное к нему линейное преобразование, определяемое тождественным равенством скалярных произведений:

$$(Ax, y) = (x, A'y), \quad \forall x, y \in E_n. \quad (75)$$

В частности,

$$(Ax, x) = (x, A'x), \quad \forall x \in E_n.$$

Преобразование называется самосопряженным, если

$$(Ax, y) = (x, Ay), \quad \forall x, y \in E_n. \quad (76)$$

Матрицы сопряженных преобразований в ортонормированном базисе связаны простым сопоставлением:

$$a_{ij}' = a_{ji}, \quad \forall i, j = 1, \dots, n,$$

Как известно, любая матрица представима в виде:

$$A = \tilde{A} + \tilde{A}', \quad (77)$$

где

$$\tilde{A} = \frac{A + A'}{2} = \tilde{A}', \quad \tilde{A} = \frac{A - A'}{2} = -\tilde{A}'. \quad (78)$$

Нетрудно видеть, что

$$B = \tilde{A} + \tilde{A}'.$$

Глава 1. §3.

- 20 -

$$0 < \tau < \tau_0 = \inf_{x \neq 0} \frac{(Bx, x)}{(Ax, x)}. \quad (87)$$

После этих замечаний перейдем к доказательству теоремы. Выразим из соотношения (69) x_k через z_k :

$$x_k = z_k + x. \quad (88)$$

и подставим в рекуррентную формулу для итерационной последовательности (65). В результате получим:

$$B \frac{z_{k+1} - z_k}{\tau} + Az_k = 0. \quad (88)$$

Отличие итерационной формулы (88) от (65) заключается в том, что она является однородной.

Матрица B – положительно определенная. Следовательно она неяврокрасивна и имеет обратную B^{-1} . С ее помощью рекуррентное соотношение (88) можно разрешить относительно z_{k+1} :

$$z_{k+1} = z_k - \tau B^{-1}Az_k = z_k - \tau \omega_k, \quad (89)$$

где

$$\omega_k = B^{-1}Az_k, \text{ так что } Az_k = B\omega_k. \quad (90)$$

Умножая обе части равенства (89) слева на матрицу A , получим еще одно рекуррентное соотношение

$$Az_{k+1} = Az_k - \tau A\omega_k. \quad (91)$$

Рассмотрим последовательность полигонтических функционалов:

$$J_k = (A\omega_k, \omega_k). \quad (92)$$

Составим аналогичные выражения для J_{k+1} преобразуем его с помощью рекуррентных формул (89) и (91):

$$J_{k+1} = (Az_k - \tau A\omega_k, z_k - \tau \omega_k) = (Az_k, z_k) - \tau (Aw_k, z_k) - \tau (Az_k, \omega_k) + \tau^2 (A\omega_k, \omega_k). \quad (93)$$

Из самосопряженности матрицы A и формулы (90) следует

$$(Aw_k, z_k) = (Az_k, \omega_k) = (B\omega_k, \omega_k).$$

В результате формула (93) принимает вид:

$$J_{k+1} = J_k - 2\tau (B\omega_k, \omega_k) + \tau^2 (A\omega_k, \omega_k) = J_k - 2\tau \left(B - \frac{\tau}{2} A \right) \omega_k, \quad (94)$$

Таким образом, последовательность функционалов J_k с учетом условия $B - \frac{\tau}{2} A > 0$ разбивает монотонно неубывающую последовательность, ограниченную сверху нулем

$$J_k \geq J_{k+1} \geq \dots \geq 0. \quad (95)$$

Поэтому она сходится. Далее, согласно лемме 3

$$\left(B - \frac{\tau}{2} A \right) \omega_k, \quad \|\omega_k\| \geq \delta \|\omega_k\|^2,$$

Глава 1. §3.

- 23 -

т. е. $\|\omega_k\| \neq 0$. Необходимость выполнения неравенства (106) для всех собственных значений μ_i для сходимости метода простой итерации доказана.

Лемма 2 определяет программу дальнейшего исследования сходимости метода простой итерации: нужно установить диапазон изменения параметра τ при котором все собственные значения удовлетворяют неравенству (106). Это легко сделать. На рис. 1 приведены графики убывающих линейных функций $\mu_i(\tau)$ (104). Все они выходят из одной точки $\tau_0 = 0$, $\mu_1 = 1$ и идут вниз из-за отрицательных коэффициентов при τ , причем быстрее всех убывает функция $\mu_1(\tau)$. Когда она принимает значение (-1) , условие (106) для нее перестает выполняться:

$$\mu_1(\tau) = 1 - \tau \lambda_1 = -1, \text{ при } \tau = \tau_0 = 2/\lambda_1.$$

Найденное значение τ_0 является границей интервала сходимости метода простой итерации

$$0 < \tau < \tau_0 = 2/\lambda_1. \quad (108)$$

Это неравенство нам уже известно. Оно было получено ранее из теоремы Самарского как достаточное условие сходимости. Дополнительный анализ на основе леммы 2 позволяет уточнить результат. Теперь мы установили, что принадлежность итерационного параметра τ интервалу (108) является необходимым и достаточным условием сходимости метода простой итерации.

Перейдем к исследованию скорости сходимости метода. Оценка погрешности (107) показывает, что она убывает на законе геометрической прогрессии со знаменателем

$$q_\tau(\tau) = \|S\| = \max_{1 \leq i \leq n} |\mu_i(\tau)|.$$

Рассмотрим рис. 2, который поможет нам провести анализ этой формулы. Он аналогичен рис. 1, только на нем приведены графики не функций $\mu_i(\tau)$, а их модулей. При малых τ все собственные значения $\mu_i(\tau)$ (104) положительны, причем наибольшим из них является $\mu_n(\tau)$, которое убывает с ростом τ с наименьшей скоростью. Однако с переходом через точку $\tau_0/2$ собственное значение $\mu_1(\tau)$, меняя знак, становится отрицательным. В результате теперь его модуль с увеличением τ не убывает, а растет и при $\tau \rightarrow \tau_0$ приближается к предельному значению – к единице.

Найдем на отрезке $[\frac{\tau_0}{2}, \tau_0]$ точку τ_* , в которой убывающая функция $\mu_1(\tau)$ сравнивается с возрастающей функцией $|\mu_1(\tau)| = -\mu_1(\tau)$. Она определяется уравнением

$$\mu_n(\tau) = 1 - \tau \lambda_n = -\mu_1(\tau) = \tau \lambda_1 - 1,$$

которое дает

$$\tau_* = \frac{2}{\lambda_1 + \lambda_n} < \tau_0. \quad (109)$$

В результате получаем:

$$\|S\| = \frac{2}{\lambda_1 + \lambda_n} < \tau_0.$$

График 1. §3.

- 24 -

График 2. §3.

- 25 -

т. е. $\|S\| \neq 0$. Необходимость выполнения неравенства (106) для всех собственных значений μ_i для сходимости метода простой итерации доказана.

Лемма 2 определяет программу дальнейшего исследования сходимости метода простой итерации: нужно установить диапазон изменения параметра τ при котором все собственные значения удовлетворяют неравенству (106). Это легко сделать. На рис. 1 приведены графики убывающих линейных функций $\mu_i(\tau)$ (104). Все они выходят из одной точки $\tau_0 = 0$, $\mu_1 = 1$ и идут вниз из-за отрицательных коэффициентов при τ , причем быстрее всех убывает функция $\mu_1(\tau)$. Когда она принимает значение (-1) , условие (106) для нее перестает выполняться:

$$\mu_1(\tau) = 1 - \tau \lambda_1 = -1, \text{ при } \tau = \tau_0 = 2/\lambda_1.$$

Найденное значение τ_0 является границей интервала сходимости метода простой итерации

$$0 < \tau < \tau_0 = 2/\lambda_1. \quad (108)$$

Это неравенство нам уже известно. Оно было получено ранее из теоремы Самарского как достаточное условие сходимости. Дополнительный анализ на основе леммы 2 позволяет уточнить результат. Теперь мы установили, что принадлежность итерационного параметра τ интервалу (108) является необходимым и достаточным условием сходимости метода простой итерации.

Перейдем к исследованию скорости сходимости метода. Оценка погрешности

$$|\mu_1(\tau)| = \|S\| = \max_{1 \leq i \leq n} |\mu_i(\tau)|$$

рассмотрим рис. 2, который поможет нам провести анализ этой формулы. Он аналогичен рис. 1, только на нем приведены графики не функций $\mu_i(\tau)$, а их модулей. При малых τ все собственные значения $\mu_i(\tau)$ (104) положительны, причем наибольшим из них является $\mu_n(\tau)$, которое убывает с ростом τ с наименьшей скоростью. Однако с переходом через точку $\tau_0/2$ собственное значение $\mu_1(\tau)$, меняя знак, становится отрицательным. В результате теперь его модуль с увеличением τ не убывает, а растет и при $\tau \rightarrow \tau_0$ приближается к предельному значению – к единице.

Найдем на отрезке $[\frac{\tau_0}{2}, \tau_0]$ точку τ_* , в которой убывающая функция $\mu_1(\tau)$ сравнивается с возрастающей функцией $|\mu_1(\tau)| = -\mu_1(\tau)$. Она определяется уравнением

$$\mu_n(\tau) = 1 - \tau \lambda_n = -\mu_1(\tau) = \tau \lambda_1 - 1,$$

которое дает

$$\tau_* = \frac{2}{\lambda_1 + \lambda_n} < \tau_0.$$

В результате получаем:

$$\|S\| = \frac{2}{\lambda_1 + \lambda_n} < \tau_0.$$

График 1. §3.

- 24 -

График 2. §3.

- 25 -

т. е. $\|S\| \neq 0$. Необходимость выполнения неравенства (106) для всех собственных значений μ_i для сходимости метода простой итерации доказана.

Лемма 2 определяет программу дальнейшего исследования сходимости метода простой итерации: нужно установить диапазон изменения параметра τ при котором все собственные значения удовлетворяют неравенству (106). Это легко сделать. На рис. 1 приведены графики убывающих линейных функций $\mu_i(\tau)$ (104). Все они выходят из одной точки $\tau_0 = 0$, $\mu_1 = 1$ и идут вниз из-за отрицательных коэффициентов при τ , причем быстрее всех убывает функция $\mu_1(\tau)$. Когда она принимает значение (-1) , условие (106) для нее перестает выполняться:

$$\mu_1(\tau) = 1 - \tau \lambda_1 = -1, \text{ при } \tau = \tau_0 = 2/\lambda_1.$$

Найденное значение τ_0 является границей интервала сходимости метода простой итерации

$$0 < \tau < \tau_0 = 2/\lambda_1. \quad (108)$$

Это неравенство нам уже известно. Оно было получено ранее из теоремы Самарского как достаточное условие сходимости. Дополнительный анализ на основе леммы 2 позволяет уточнить результат. Теперь мы установили, что принадлежность итерационного параметра τ интервалу (108) является необходимым и достаточным условием сходимости метода простой итерации.

Перейдем к исследованию скорости сходимости метода. Оценка погрешности

$$|\mu_1(\tau)| = \|S\| = \max_{1 \leq i \leq n} |\mu_i(\tau)|$$

рассмотрим рис. 2, который поможет нам провести анализ этой формулы. Он аналогичен рис. 1, только на нем приведены графики не функций $\mu_i(\tau)$, а их модулей. При малых τ все собственные значения $\mu_i(\tau)$ (104) положительны, причем наибольшим из них является $\mu_n(\tau)$, которое убывает с ростом τ с наименьшей скоростью. Однако с переходом через точку $\tau_0/2$ собственное значение $\mu_1(\tau)$, меняя знак, становится отрицательным. В результате теперь его модуль с увеличением τ не убывает, а растет и при $\tau \rightarrow \tau_0$ приближается к предельному значению – к единице.

Найдем на отрезке $[\frac{\tau_0}{2}, \tau_0]$ точку τ_* , в которой убывающая функция $\mu_1(\tau)$ сравнивается с возрастающей функцией $|\mu_1(\tau)| = -\mu_1(\tau)$. Она определяется уравнением

$$\mu_n(\tau) = 1 - \tau \lambda_n = -\mu_1(\tau) = \tau \lambda_1 - 1,$$

которое дает

$$\tau_* = \frac{2}{\lambda_1 + \lambda_n} < \tau_0.$$

В результате получаем:

$$\|S\| = \frac{2}{\lambda_1 + \lambda_n} < \tau_0.$$

График 1. §3.

- 24 -

График 2. §3.

- 25 -

т. е. $\|S\| \neq 0$. Необходимость выполнения неравенства (106) для всех собственных значений μ_i для сходимости метода простой итерации доказана.

Лемма 2 определяет программу дальнейшего исследования сходимости метода простой итерации: нужно установить диапазон изменения параметра τ при котором все собственные значения удовлетворяют неравенству (106). Это легко сделать. На рис. 1 приведены графики убывающих линейных функций $\mu_i(\tau)$ (104). Все они выходят из одной точки $\tau_0 = 0$, $\mu_1 = 1$ и идут вниз из-за отрицательных коэффициентов при τ , причем быстрее всех убывает функция $\mu_1(\tau)$. Когда она принимает значение (-1) , условие (106) для нее перестает выполняться:

$$\mu_1(\tau) = 1 - \tau \lambda_1 = -1, \text{ при } \tau = \tau_0 = 2/\lambda_1.$$

Найденное значение τ_0 является границей интервала сходимости метода простой итерации

$$0 < \tau < \tau_0 = 2/\lambda_1. \quad (108)$$

Это неравенство нам уже известно. Оно было получено ранее из теоремы Самарского как достаточное условие сходимости. Дополнительный анализ на основе леммы 2 позволяет уточнить результат. Теперь мы установили, что принадлежность итерационного параметра τ интервалу (108) является необходимым и достаточным условием сходимости метода простой итерации.

Перейдем к исследованию скорости сходимости метода. Оценка погрешности

$$|\mu_1(\tau)| = \|S\| = \max_{1 \leq i \leq n} |\mu_i(\tau)|$$

рассмотрим рис. 2, который поможет нам провести анализ этой формулы. Он аналогичен рис. 1, только на нем приведены графики не функций $\mu_i(\tau)$, а их модулей. При малых τ все собственные значения $\mu_i(\tau)$ (104) положительны, причем наибольшим из них является $\mu_n(\tau)$, которое убывает с ростом τ с наименьшей скоростью. Однако с переходом через точку $\tau_0/2$ собственное значение $\mu_1(\tau)$, меняя знак, становится отрицательным. В результате теперь его модуль с увеличением τ не убывает, а растет и при $\tau \rightarrow \tau_0$ приближается к предельному значению – к единице.

Найдем на отрезке $[\frac{\tau_0}{2}, \tau_0]$ точку τ_* , в которой убывающая функция $\mu_1(\tau)$ сравнивается с возрастающей функцией $|\mu_1(\tau)| = -\mu_1(\tau)$. Она определяется уравнением

$$\mu_n(\tau) = 1 - \tau \lambda_n = -\mu_1(\tau) = \tau \lambda_1 - 1,$$

которое дает

$$\tau_* = \frac{2}{\lambda_1 + \lambda_n} < \tau_0.$$

В результате получаем:

$$\|S\| = \frac{2}{\lambda_1 + \lambda_n} < \tau_0.$$

График 1. §3.

- 24 -

График 2. §3.

- 25 -

т. е. $\|S\| \neq 0$. Необходимость выполнения неравенства (106) для всех собственных значений μ_i для сходимости метода простой итерации доказана.

Лемма 2 определяет программу дальнейшего исследования сходимости метода простой итерации: нужно установить диапазон изменения параметра τ при котором все собственные значения удовлетворяют неравенству (106). Это легко сделать. На рис. 1 приведены графики убывающих линейных функций $\mu_i(\tau)$ (104). Все они выходят из одной точки $\tau_0 = 0$, $\mu_1 = 1$ и идут вниз из-за отрицательных коэффициентов при τ , причем быстрее всех убывает функция $\mu_1(\tau)$. Когда она принимает значение (-1) , условие (106) для нее перестает выполняться:

$$\mu_1(\tau) = 1 - \tau \lambda_1 = -1, \text{ при } \tau = \tau_0 = 2/\lambda_1.$$

Найденное значение τ_0 является границей интервала сходимости метода простой итерации

$$0 < \tau < \tau_0 = 2/\lambda_1. \quad (108)$$

Это неравенство нам уже известно. Оно было получено ранее из теоремы Самарского как достаточное условие сходимости. Дополнительный анализ на основе леммы 2 позволяет уточнить результат. Теперь мы установили, что принадлежность итерационного параметра τ интервалу (108) является необходимым и достаточным условием сходимости метода простой итерации.

Перейдем к исследованию скорости сходимости метода. Оценка погрешности

$$D_g = a_g \delta_g = \begin{cases} 0, & i \neq j \\ a_{ii}, & i = j \end{cases}$$

 T_H - нижняя треугольная матрица

$$(T_H)_i = \begin{cases} a_{ii}, & i > j \\ 0, & i \leq j \end{cases}$$

 T_B - верхняя треугольная матрица

$$(T_B)_i = \begin{cases} 0, & i \geq j \\ a_{ii}, & i < j \end{cases}$$

В классическом методе Зейделя, записанном в канонической форме, полагают

$$B = D + T_H, \quad (115)$$

 $\tau = 1$.

В результате формула (65) принимает вид:

$$(D + T_H)(\mathbf{x}_{k+1} - \mathbf{x}_k) + A\mathbf{x}_k = \mathbf{f},$$

или

$$(D + T_H)\mathbf{x}_{k+1} + T_B\mathbf{x}_k = \mathbf{f}. \quad (116)$$

Перейдем от векторной формы записи рекуррентной формулы (116) к построчной:

$$\begin{aligned} a_{11}x_1^{k+1} + a_{12}x_2^{k+1} + a_{13}x_3^{k+1} + \cdots + a_{1n}x_n^{k+1} &= f_1 \\ a_{21}x_1^{k+1} + a_{22}x_2^{k+1} + a_{23}x_3^{k+1} + \cdots + a_{2n}x_n^{k+1} &= f_2 \\ \vdots & \vdots \\ a_{n1}x_1^{k+1} + a_{n2}x_2^{k+1} + a_{n3}x_3^{k+1} + \cdots + a_{nn}x_n^{k+1} &= f_n. \end{aligned} \quad (117)$$

Уравнения (117) позволяют последовательно рассчитать компоненты вектора $(k+1)$ -ой итерации подобно тому, как это делалось во время обратного хода в методе Гаусса:

$$x_i^{k+1} = \frac{1}{a_{ii}} \left[f_i - \sum_{j=1}^{i-1} a_{ij}x_j^{k+1} - \sum_{j=i+1}^n a_{ij}x_j^k \right], \quad i = 1, \dots, n. \quad (118)$$

Формула (118) предполагает, что $a_{ii} \neq 0$, $1 \leq i \leq n$. Если матрица A удовлетворяет условию теоремы Самарского (84), $A - A' > 0$, то, согласно неравенству (81), все ее диагональные элементы должны быть строго положительными и, тем самым, не могут обращаться в нуль.Алгоритм в методе Зейделя прост и удобен для вычислений. Он не требует никаких действий с матрицей A . Ранее вычисленные на текущей итерации компоненты x_j^{k+1} ($j < i$) сразу же участвуют в расчетах наряду с компонентами x_j^k ($j > i$) и, таким образом, не требуют дополнительного резерва памяти, что существенно при решении больших систем.Сходимость метода Зейделя в случае, когда матрица A удовлетворяет условию теоремы Самарского, т.е. является самосопряженной и положительно определенной, будет доказана в следующем разделе. К этому утверждению добавим без доказательства еще один результат: метод Зейделя сходится для любой системы (62), в которой матрица A обладает свойством диагонального преобладания.

$$\frac{3}{4}D + T_H = \begin{bmatrix} 3/4 & 0 \\ 1 & 3/2 \end{bmatrix}, \quad \frac{1}{4}D + T_B = \begin{bmatrix} 1/4 & 1 \\ 0 & 1/2 \end{bmatrix}.$$

С их помощью рекуррентное соотношение (123), записанное покомпонентно, принимает вид:

$$\begin{aligned} \frac{3}{4}x_1^{k+1} + \frac{1}{4}x_2^k + x_2^{k+1} &= 0, \\ x_1^{k+1} + \frac{3}{2}x_2^{k+1} + \frac{1}{2}x_2^k &= 1. \end{aligned}$$

Выражая из первого соотношения x_1^{k+1} , из второго x_2^{k+1} , получим окончательные расчетные формулы для компонент очередной итерации:

$$\begin{aligned} x_1^{k+1} &= -\frac{3}{3}x_2^k - \frac{4}{3}x_2^k, \\ x_2^{k+1} &= \frac{2}{3}x_1^{k+1} - \frac{1}{3}x_2^k. \end{aligned}$$

Примем, как и в предыдущих случаях, за начальное приближение нулевой вектор и сделаем три итерации. При этом для каждой из них подсчитаем невязку (71), позволяющую следить за сходимостью процесса:

$$\begin{aligned} x_1 &= \begin{pmatrix} 0 \\ 2 \\ 3 \end{pmatrix}, \quad \Psi_1 = \begin{pmatrix} 2/3 \\ 1/3 \\ 1/3 \end{pmatrix}, \quad \|\Psi_1\| = \sqrt{\frac{13}{3}} \approx 0.745, \\ x_2 &= \begin{pmatrix} 8/27 \\ 1/9 \\ 27/27 \end{pmatrix}, \quad \Psi_2 = \begin{pmatrix} 4/27 \\ 5/27 \\ 22/27 \end{pmatrix}, \quad \|\Psi_2\| = \sqrt{\frac{41}{27}} \approx 0.237, \\ x_3 &= \begin{pmatrix} -88/256 \\ 8/243 \\ 243/243 \end{pmatrix}, \quad \Psi_3 = \begin{pmatrix} -8/243 \\ 5/243 \\ 243/243 \end{pmatrix}, \quad \|\Psi_3\| = \sqrt{\frac{89}{243}} \approx 0.039. \end{aligned}$$

Поведение невязок, а также сравнение членов итерационной последовательности x_i с точным решением системы $x = [-1, 1]$ показывают сходимость процесса, более быстрое, чем в методе Зейделя. Выбранное значение параметра $\omega = 4/3$ оказалось близким к оптимальному $\omega = \omega_*$.Согласно рассмотренной выше общей схеме построения интерполяирующей функции, следует потребовать, чтобы коэффициенты c_i с учетом (7) удовлетворяли системе линейных уравнений:

$$\sum_{i=0}^n c_i x_i^j = f(x_i), \quad j = 0, 1, \dots, n. \quad (8)$$

Определителем этой системы является определитель Ван-дер-Монда:

$$\Delta = \begin{vmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ & 1 & x_1 & x_1^2 & \cdots & x_1^n \\ & & \vdots & \vdots & \ddots & \vdots \\ & & x_n & x_n^2 & \cdots & x_n^n \end{vmatrix} = \prod_{i>j} (x_i - x_j).$$

В нашем случае этот определитель отличен от нуля, поскольку, согласно (1), все узлы интерполяции различны между собой. Итак, интерполяирование с помощью полиномов при сделанных в начале главы предположениях всегда осуществляется и притом единственным образом.

Задача 1.

Построить линейный полином

$$P_1(x) = c_0 + c_1 x$$

по заданным узлам интерполяции $x_0 < x_1$ и соответствующим им значениям функции

$$y_0 = f(x_0) \text{ и } y_1 = f(x_1).$$

Линейная система уравнений для определения c_0 и c_1 в данном случае имеет вид:

$$c_0 + c_1 x_0 = f(x_0),$$

$$c_0 + c_1 x_1 = f(x_1).$$

Определитель этой системы равен $\Delta = x_1 - x_0 > 0$. Решив систему, получим:

$$c_0 = \frac{x_1 f(x_0) - x_0 f(x_1)}{x_1 - x_0}, \quad c_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0}.$$

Следовательно,

$$P_1(x) = \frac{x_1 f(x_0) - x_0 f(x_1)}{x_1 - x_0} + \frac{f(x_1) - f(x_0)}{x_1 - x_0} x. \quad (9)$$

Перепишем этот полином в несколько другой форме, выделив $f(x_0)$ и $f(x_1)$ в качестве множителей

$$P_1(x) = f(x_0) \frac{x - x_1}{x_0 - x_1} + f(x_1) \frac{x - x_0}{x_1 - x_0}. \quad (10)$$

Геометрический образ интерполяирующей функции $P_1(x)$ - прямая, проходящая на плоскости (x, y) через точки с координатами (x_0, y_0) и (x_1, y_1) . Уравнение этой прямой, наряду с (9) и (10), можно переписать в виде:

$$y = P_1(x) = f(x_0) + \frac{f(x_1) - f(x_0)}{x_1 - x_0} (x - x_0). \quad (11)$$

Задача 3.

Рассмотреть систему (112) (см. задачу 2) и построить для нее приближенное решение с помощью метода Зейделя.

В рассматриваемом случае рекуррентные формулы (118) для построения $(k+1)$ -ой итерации по k -ой итерации принимают вид:

$$\begin{aligned} x_1^{k+1} &= -x_2^k \\ x_2^{k+1} &= \frac{1}{2}(1 - x_1^{k+1}). \end{aligned} \quad (119)$$

Принимая, как при решении задачи 2, за начальное приближение нулевой вектор, подсчитаем по формуле (119) несколько первых итераций, сопроводив этот процесс подсчетом невязки:

$$\begin{aligned} x_1 &= \begin{pmatrix} 0 \\ 1/2 \end{pmatrix}, \quad \Psi_1 = \begin{pmatrix} 1/2 \\ 0 \end{pmatrix}, \quad \|\Psi_1\| = \frac{1}{2}, \\ x_2 &= \begin{pmatrix} -1/2 \\ 1/4 \end{pmatrix}, \quad \Psi_2 = \begin{pmatrix} 1/4 \\ 0 \end{pmatrix}, \quad \|\Psi_2\| = \frac{1}{4}, \\ x_3 &= \begin{pmatrix} -3/8 \\ 7/16 \end{pmatrix}, \quad \Psi_3 = \begin{pmatrix} 1/8 \\ 0 \end{pmatrix}, \quad \|\Psi_3\| = \frac{1}{8}. \end{aligned}$$

Обсудим полученные результаты. Начнем с невязки. Ее вторая компонента все время остается равной нулю, поскольку второе уравнение системы на каждой итерации выполняется, как видно из (119), точно. Первые компоненты невязки и звездочки убывают по закону геометрической прогрессии с знаменателем 1/2, т.е. гораздо быстрее, чем в методе простой итерации. Хорошая сходимость процесса видна также из прямого сравнения чисел итерационной последовательности x_i с точным решением системы $x = [-1, 1]$.**3.6. Метод верхней релаксации**Модифицируем метод Зейделя. С этой целью введем параметр ω и записем рекуррентное соотношение (65) в виде

$$(D + \omega T_H) \frac{(x_{i-1} - x_i)}{\omega} + Ax_i = \mathbf{f}. \quad (120)$$

В данном случае

$$B = D + \omega T_H, \quad \tau = \omega > 0. \quad (121)$$

При $\omega = 1$ мы возвращаемся к методу Зейделя.

Соотношение (120) можно придать вид

$$\left(B - \frac{\tau}{2} A \right) \frac{(x_{i-1} - x_i)}{\omega} + Ax_i = \mathbf{f}. \quad (122)$$

Такая форма записи показывает, что параметр ω влияет на диагональ матрицы B .Для построения алгоритма вычисления очередной итерации нужно разделить влево части рекуррентной формулы (122) члены, содержащие x_{i-1} и x_i , и придать ей форму, аналогичную (116):

$$\left(\frac{1}{\omega} D + T_H \right) \frac{(x_{i-1} - x_i)}{\omega} + Ax_i = \mathbf{f}. \quad (123)$$

Сходимость метода Зейделя в случае, когда матрица A удовлетворяет условию теоремы Самарского, т.е. является самосопряженной и положительно определенной, будет доказана в следующем разделе. К этому утверждению добавим без доказательства еще один результат: метод Зейделя сходится для любой системы (62), в которой матрица A обладает свойством диагонального преобладания.

$$\frac{3}{4}D + T_H = \begin{bmatrix} 3/4 & 0 \\ 1 & 3/2 \end{bmatrix}, \quad \frac{1}{4}D + T_B = \begin{bmatrix} 1/4 & 1 \\ 0 & 1/2 \end{bmatrix}.$$

С их помощью рекуррентное соотношение (123), записанное покомпонентно, принимает вид:

$$\begin{aligned} \frac{3}{4}x_1^{k+1} + \frac{1}{4}x_2^k + x_2^{k+1} &= 0, \\ x_1^{k+1} + \frac{3}{2}x_2^{k+1} + \frac{1}{2}x_2^k &= 1. \end{aligned}$$

Выражая из первого соотношения x_1^{k+1} , из второго x_2^{k+1} , получим окончательные расчетные формулы для компонент очередной итерации:

$$\begin{aligned} x_1^{k+1} &= -\frac{3}{3}x_2^k - \frac{4}{3}x_2^k, \\ x_2^{k+1} &= \frac{2}{3}x_1^{k+1} - \frac{1}{3}x_2^k. \end{aligned}$$

Примем, как и в предыдущих случаях, за начальное приближение нулевой вектор и сделаем три итерации. При этом для каждой из них подсчитаем невязку (71), позволяющую следить за сходимостью процесса

$$\begin{aligned} x_1 &= \begin{pmatrix} 0 \\ 2 \\ 3 \end{pmatrix}, \quad \Psi_1 = \begin{pmatrix} 2/3 \\ 1/3 \\ 1/3 \end{pmatrix}, \quad \|\Psi_1\| = \sqrt{\frac{13}{3}} \approx 0.745, \\ x_2 &= \begin{pmatrix} 8/27 \\ 1/9 \\ 27/27 \end{pmatrix}, \quad \Psi_2 = \begin{pmatrix} 4/27 \\ 5/27 \\ 22/27 \end{pmatrix}, \quad \|\Psi_2\| = \sqrt{\frac{41}{27}} \approx 0.237, \\ x_3 &= \begin{pmatrix} -88/256 \\ 8/243 \\ 243/243 \end{pmatrix}, \quad \Psi_3 = \begin{pmatrix} -8/243 \\ 5/243 \\ 243/243 \end{pmatrix}, \quad \|\Psi_3\| = \sqrt{\frac{89}{243}} \approx 0.039. \end{aligned}$$

Поведение невязок, а также сравнение членов итерационной последовательности x_i с точным решением системы $x = [-1, 1]$ показывают сходимость процесса, более быстрое, чем в методе Зейделя. Выбранное значение параметра $\omega = 4/3$ оказалось близким к оптимальному $\omega = \omega_*$.**Глава 2. ПРИБЛИЖЕНИЕ ФУНКЦИЙ.**Пусть на отрезке $[a, b]$ определена некоторая функция $y = f(x)$, однако полная информация о ней недоступна. Известны лишь ее значения в конечном числе точек x_1, x_2, \dots, x_n , этого отрезка, которые мы будем считать замкнутыми в порядке возрастания:

$$a \leq x_0 < x_1 < \dots < x_{n-1} < x_n \leq b. \quad (1)$$

Требуется по известным значениям

$$y_i = f(x_i), \quad i = 0, 1, \dots, n \quad (2)$$

«восстановить», хотя бы приближенно, исходную функцию $y = f(x)$, то есть построить на отрезке $[a, b]$ функцию $F(x)$, достаточно близкую к $f(x)$. Функцию $F(x)$ принято называть интерполярующей функцией, точки $x = x_0, x = x_1, \dots, x_n$ - узлами интерполяции.Подобные задачи часто возникают на практике, например, при обработке экспериментальных данных, когда значения переменной y , зависящей от x , измеряются в конечном числе точек x_i , $y_i = f(x_i)$, $i = 0, 1, \dots, n$ или при работе с табличными функциями, если требуется вычислить $y = f(x)$, при значениях аргумента, не совпадающему с одним из технических x_i .

Поставленный выше в общей форме вопрос о построении функций является достаточно сложным. Существует не один подход к его решению. Мы ограничимся изложением трех наиболее распространенных методов.

§1. Интерполяция.**1.1. Классическая постановка задачи интерполяирования.**Выберем некоторую систему функций $\varphi_0(x), \varphi_1(x), \dots, \varphi_n(x)$, заданную на отрезке $[a, b]$, и будем строить $F(x)$ как их линейную комбинацию:

$$F(x) = \sum_{i=0}^n c_i \varphi_i(x), \quad (3)$$

где числовые коэффициенты c_i , $i = 0, 1, \dots, n$ подлежат определению, согласно условиям:

$$F(x_i) = f(x_i), \quad i = 0, 1, \dots, n. \quad (4)$$

Равенства (4) представляют собой систему линейных алгебраических уравнений относительно коэффициентов c_i :

$$\sum_{i=0}^n c_i \varphi_i(x_j) = f(x_j), \quad j = 0, 1, \dots, n, \quad (5)$$

или в развернутом виде:

$$Q_{n,0}(x) = \frac{(x - x_1)(x - x_2) \cdots (x - x_n)}{(x_0 - x_1)(x_0 - x_2) \cdots (x_0 - x_n)},$$

$$Q_{n,1}(x) = \frac{(x - x_0)(x - x_2) \cdots (x - x_n)}{(x_1 - x_0)(x_1 - x_2) \cdots (x_1 - x_n)},$$

$$Q_{n,n}(x) = \frac{(x - x_0)(x - x_1) \cdots (x - x_{n-1})}{(x_n - x_0)(x_n - x_1) \cdots (x_n - x_{n-1})},$$

Иногда нам будет удобно записывать $Q_{n,i}(x)$ в виде:

$$Q_{n,i}(x) = \frac{(x - x_0) \cdots [i] \cdots (x - x_n)}{(x_i - x_0) \cdots [i] \cdots (x_i - x_n)}.$$

Из выражения (12) и формулы (13) очевидно, что построенный полином $P_1(x)$ действительно является интерполяционным полиномом для функции $y = f(x)$ на сетке с узлами x_0, x_1, \dots, x_n . Его принято называть интерполяционным полиномом в форме Лагранжа.Из этого следует, что возможны и другие эквивалентные представления интерполяционного полинома $P_1(x)$. Одним из них мы познакомимся в следующем разделе.

В заключение отметим, что из трех различных представлений интерполяционного полинома первой степени (9)-(11) формула (10) дает его запись в форме Лагранжа.

Задача 2.

Из данного примера видно, что всегда существуют различные эквивалентные между собой формы записи интерполяционного полинома, удобные в различных ситуациях.

1.3. Построение интерполяционного полинома в форме Лагранжа.Интерполяционный полином первой степени (9) мы построили решая на прямую систему двух уравнений с двумя неизвестными - коэффициентами c_0 и c_1 . Однако решить такую же систему (8) при произвольном n технически очень сложно. Попробуем это сделать, это с помощью специальных методов, учитывающих особенности рассматриваемой задачи. Одни из таких методов, принадлежащие Лагранжу, мы рассмотрим в этом разделе.Представим искомый полином $P_n(x)$ в виде:

$$P_n(x) = \sum_{i=0}^n f(x_i) Q_{n,i}(x), \quad (12)$$

где $Q_{n,i}$ полиномы степени n , «ориентированные» на точке x_i в том смысле, что

$$Q_{n,i}(x_i) = \begin{cases} 0, & x_i = x_j, \quad \forall j \neq i, \\ 1, & x_i = x_i. \end{cases} \quad (13)$$

Такие полиномы легко построить:

$$Q_{n,i}(x) = \prod_{j=0, j \neq i}^n \frac{(x - x_j)}{(x_i - x_j)}, \quad (14)$$

или в развернутом виде:

$$Q_{n,i}(x) = \frac{(x - x_0) \cdots [i] \cdots (x - x_n)}{(x_i - x_0) \cdots [i] \cdots (x_i - x_n)}.$$

Из выражения (12) и формулы (13) очевидно, что построенный полином $P_n(x)$ действительно является интерполяционным полиномом для функции $y = f(x)$ наплоскости (x, y) через точки $(x_0, y_0), \dots, (x_n, y_n)$. Выведем с помощью этого

результата формулу (12).

$$P_n(x) = \sum_{i=0}^n f(x_i) Q_{n,i}(x), \quad (15)$$

где $P_n(x)$ - полином Лагранжа степени n , соответствующие узлыимеют степень n .</

где

$$\omega_{x_i} = (x_i - x_0) \dots (x_i - x_{i-1}) (x_i - x_{i+1}) \dots (x_i - x_n). \quad (22)$$

При этом

$$A_0 = f(x_0). \quad (23)$$

Формулы (19) и (21) позволяют написать рекуррентное соотношение для полинома $P_n(x)$:

$$P_n(x) = P_{n-1}(x) + A_n(x - x_0) \dots (x - x_{n-1}). \quad (24)$$

Выражая аналогичным образом по индукции $P_n(x)$ через $P_{n-2}(x)$, $P_{n-1}(x)$ через $P_{n-3}(x)$ и т. д., получим окончательную формулу для полинома $P_n(x)$:

$$P_n(x) = A_0 + A_1(x - x_0) + A_2(x - x_0)(x - x_1) + \dots + A_{n-1}(x - x_0) \dots (x - x_{n-2})(x - x_{n-1}). \quad (25)$$

Представление (25) удобно для вычисления, поскольку увеличение n на единицу требует только добавления к «старому» многочлену одного дополнительного слагаемого. Такое представление интерполяционного полинома $P_n(x)$ называют интерполяционным полиномом в форме Ньютона.

Из трех эквивалентных представлений интерполяционного полинома первой степени (9) – (11) формула (11) дает его запись в форме Ньютона.

Задача 3.Написать интерполяционный полином второй степени в форме Ньютона для функции $y = \sin x$ по ее значениям в трех точках: $x_0 = 0$, $x_1 = \pi/6$, $x_2 = \pi/2$ (см. задачу 2).

Согласно формуле (25)

$$P_2(x) = A_0 + A_1x + A_2x^2 \left(x - \frac{\pi}{6} \right). \quad (26)$$

Коэффициенты в этом разложении вычисляются по формулам (21) и (23):

$$A_0 = 0, \quad A_1 = \frac{1}{2} \left(\frac{\pi}{6} \right), \quad A_2 = -\frac{18}{2\pi^2} + \frac{6}{\pi^3} = -\frac{3}{\pi^2}. \quad (27)$$

Подставляя найденные значения коэффициентов в формулу (26), получим

$$P_2(x) = \frac{3}{\pi} x - \frac{3x^2}{\pi^2} \left(x - \frac{\pi}{6} \right) = \frac{x}{\pi^2} \left(\frac{7\pi}{2} - 3x \right). \quad (28)$$

Первоначальные выражения для интерполяционного полинома в форме Лагранжа и Ньютона различны, но окончательные ответы, естественно, совпадают.

5.1. Погрешность интерполяирования.Поставим вопрос о том, насколько хорошо интерполяционный полином $P_n(x)$ приближает функцию $f(x)$ на отрезке $[a, b]$, то есть попытаемся оценить погрешность (остаточный член)

$$R_n(x) = f(x) - P_n(x), \quad x \in [a, b]. \quad (29)$$

$$R_2\left(\frac{\pi}{4}\right) \leq \frac{\pi^3}{1152} < 0.027.$$

Эта оценка согласуется с величиной погрешности (17), вычисленной «в лоб».

1.6. О сходимости интерполяционного процесса.Поставим вопрос, будут ли сходиться интерполяционные полиномы $P_n(x)$ к интерполяируемой функции $f(x)$ на отрезке $[a, b]$ при неограниченном возрастании числа узлов n .Упорядоченное множество точек x_i , $i = 0, 1, \dots, n$ назовем сеткой на отрезке $[a, b]$ и обозначим для краткости Ω_n . Рассмотрим последовательность сеток с возрастающим числом узлов:

$$\Omega_0 = \{x_0^{(0)}\}, \quad \Omega_1 = \{x_0^{(1)}, x_1^{(1)}\}, \dots, \Omega_n = \{x_0^{(n)}, x_1^{(n)}, \dots, x_n^{(n)}\}, \dots$$

и отвечающую ей последовательность интерполяционных полиномов $P_n(x)$, построенных для фиксированной непрерывной на отрезке $[a, b]$ функции $f(x)$.Интерполяционный процесс для функции сходится в точке $x \in [a, b]$, если существует предел

$$\lim_{n \rightarrow \infty} P_n(x) = f(x).$$

Наряду с обычной сходимостью часто рассматривается сходимость в различных нормах. Так, равномерная сходимость на отрезке $[a, b]$ означает, что

$$\max_{x \in [a, b]} |f(x) - P_n(x)| \rightarrow 0 \text{ при } n \rightarrow \infty.$$

Сходимость или расходимость интерполяционного процесса зависит как от выбора последовательности сеток, так и от гладкости функции $f(x)$. Если $f(x)$ – цепь аналитической функции, то при произвольном расположении узлов на отрезке $[a, b]$ интерполяционный многочлен $P_n(x)$ равномерно сходится к $f(x)$ при $n \rightarrow \infty$.Положение резко меняется, если производные функции разрывны или не существуют в отдельных точках. Например для функции $f(x) = |x|$ на отрезке $[-1, 1]$, покрытом равномерной сеткой узлов, значения $P_n(x)$ между узлами интерполяции неограниченно возрастают при $n \rightarrow \infty$. Вместе с тем, для заданной непрерывной функции $f(x)$ за счет выбора сеток можно добиться сходимости и притом равномерной на $[a, b]$. Однако построение таких сеток довольно сложно и, главное, такие сетки «индивидуальны» для каждой конкретной функции.Если заметить дополнительно, что объем вычислений при построении интерполяционного полинома быстро нарастает с ростом n , то становится понятно, что на практике вычислители избегают пользоваться интерполяционными полиномами высокой степени. Вместо этого, в случае необходимости, при больших значениях n используется кусочно-полиномиальная интерполяция, которую мы обсудим в следующем параграфе.

$$H'_{n+1}(x_i) = f_i \left(\frac{1}{x_i - x_0} + \dots + \frac{1}{x_i - x_{i-1}} + \frac{1}{x_i - x_{i+1}} + \dots + \frac{1}{x_i - x_n} \right) + \alpha_j. \quad (45)$$

Для соблюдения требования $H'_{n+1}(x_i) = f'_i$ следует положить

$$\alpha_j = f'_i - f_i A_j,$$

где для краткости обозначено

$$A_j = \frac{1}{x_j - x_0} + \dots + \frac{1}{x_j - x_{j-1}} + \frac{1}{x_j - x_{j+1}} + \dots + \frac{1}{x_j - x_n}. \quad (46)$$

Итак:

$$H'_{n+1}(x) = \sum_{k=0}^n f_k \frac{(x - x_0) \dots (x - x_k) \dots (x - x_n)}{(x_k - x_0) \dots (x_k - x_{k-1}) \dots (x_k - x_{k+1}) \dots (x_k - x_n)} \left(\frac{x - x_i}{x_i - x_j} \right) + \left[f_j + (f'_i - f_i A_j) (x - x_i) \right] \frac{(x - x_0) \dots (x - x_{j-1}) \dots (x - x_{j+1}) \dots (x - x_n)}{(x_i - x_0) \dots (x_i - x_{j-1}) \dots (x_i - x_{j+1}) \dots (x_i - x_n)}. \quad (47)$$

Пример 2.Построить интерполяционный полином Эрмита для функции $f(x)$ в случае, когда во всех узлах интерполяции x_k , $k = 0, 1, \dots, m$ заданы значения функции $f(x_i) = f_i$ и ее первой производной $f'(x_i) = f'_i$.В данном случае $N_n = 2$, $k = 0, 1, \dots, m$, так что степень полинома $H_n(x)$ равна $2m+1$.

Запишем исходный полином в виде:

$$H_{2m+1}(x) = f_i \left(\frac{2}{x_i - x_0} + \dots + \frac{2}{x_i - x_{i-1}} + \frac{2}{x_i - x_{i+1}} + \dots + \frac{2}{x_i - x_n} \right) + \alpha_i = f'_i. \quad (48)$$

Представление (48) удобно тем, что автоматически выполняются условия

$$H_{2m+1}(x_i) = f_i.$$

При вычислении производной полинома (48) в узле $x = x_k$ следует учесть, что все слагаемые суммы, кроме слагаемого, отвечающего самому узлу x_k , дают нулевой вклад в производную в этой точке, поэтому

$$H'_{2m+1}(x_k) = f_i \left(\frac{2}{x_k - x_0} + \dots + \frac{2}{x_k - x_{k-1}} + \frac{2}{x_k - x_{k+1}} + \dots + \frac{2}{x_k - x_n} \right) + \alpha_k = f'_i.$$

Отсюда

$$\alpha_k = f'_i - 2f_i A_k,$$

где, числа A_k определяются формулой (46). Таким образом, решением данной задачи является полиномом Эрмита

$$H_{2m+1}(x) = \sum_{k=0}^m f_k + (f'_i - 2f_i A_k) (x - x_i) \frac{(x - x_0)^2 \dots (x - x_{i-1})^2 (x - x_{i+1})^2 \dots (x - x_n)^2}{(x_i - x_0)^2 \dots (x_i - x_{i-1})^2 (x_i - x_{i+1})^2 \dots (x_i - x_n)^2}. \quad (49)$$

Задача 5

Сразу же отметим, что по определению интерполяционного полинома

$$R_n(x_i) = 0 \text{ при } i = 0, 1, \dots, n, \quad (30)$$

поскольку речь идет об оценке $R_n(x)$ при значениях $x = x_i$.Для того, чтобы это сделать, следует ввести дополнительно предположение о гладкости функции $f(x)$. Предположим, что $f(x)$ имеет $(n+1)$ непрерывную производную на отрезке $[a, b]$.В силу (30) $R_n(x)$ можно представить в виде:

$$R_n(x) = \omega_{n+1}(x) P_n(x), \quad (31)$$

где $\omega_{n+1}(x)$ – полином степени $(n+1)$:

$$\omega_{n+1}(x) = (x - x_0)(x - x_1) \dots (x - x_n). \quad (32)$$

Зафиксируем произвольное значение $x \in [a, b]$ и рассмотрим вспомогательную функцию от переменной t :

$$g(t) = f(t) - P_n(t) - \omega_{n+1}(t) P_n(x),$$

заданную на отрезке $[a, b]$ и содержащую переменную x в качестве параметра. В силу определения функция $g(t)$ обязана обращаться в нуль в узлахинтерполяции при $t = x_i$ и кроме того при $t = x$, т. е. как функция аргумента t онаимеет нулевую производную на отрезке $[a, b]$.

$$g(x) = 0, \quad i = 0, 1, \dots, n, \quad g(x) = 0. \quad (33)$$

Если $x \in [x_0, x_n]$, то все ее нули также лежат на отрезке $[x_0, x_n]$. Если $x < x_0$, то эти нули, вообще говоря, принадлежат отрезку $[x_0, x_n]$, а если $x > x_n$, то они находятся на отрезке $[x_0, x_n]$. Объединяя эти три случая, получаем, что указанные нули функции $g(t)$ припадлежат отрезку $[a, b]$, где $a = \min(x_0, x_n) \geq a$, $b = \max(x_n, x_0) \leq b$.

$$g(t) = 0 \text{ на } [a, b], \quad t \in [a, b]. \quad (34)$$

Согласно известной теореме Ролля можем утверждать, что производная $g'(t)$ имеет по крайней мере $(n+1)$ нуль на отрезке $[a, b]$ (эти нули перемежаются с нулями самой функции $g(t)$). Повторяя это рассуждение, заключаем, что $g^{(i)}(t)$ имеет по крайней мере n нулей на отрезке $[a, b]$, $g^{(n)}(t) - (n-1)$ нуль и, наконец, $g^{(n+1)}(t)$ обращается хотя бы один раз в нуль в некоторой точке $t = \xi \in [a, b]$, то есть

$$g^{(n+1)}(\xi) = f^{(n+1)}(\xi) - P_n^{(n+1)}(\xi) - (n+1) P_n(x) = 0.$$

Учитывая, что $(n+1)$ производная полинома степени n тождественно равна нулю, получаем, что

$$r_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{n+1}(x). \quad (35)$$

Формула (35) не позволяет вычислить погрешность, поскольку точное значение аргумента ξ нам известно. Однако с ее помощью погрешность можно оценить:

$$|R_n(x)| \leq \frac{M_{n+1}}{(n+1)!} |\omega_{n+1}(x)|. \quad (36)$$

где

$$M_{n+1} = \max_{x \in [a, b]} |f^{(n+1)}(x)| \leq \max_{x \in [a, b]} |P^{(n+1)}(x)|. \quad (37)$$

Обсудим роль полинома $\omega_{n+1}(x)$ (32) в оценке (36). На отрезке $[x_0, x_n]$ он имеет $(n+1)$ нуль, а его значения между этими нулями сравнительно невелики, но, когда точка x выходит за пределы отрезка $[x_0, x_n]$ и удаляется от точки x_0 влево или от точки x_n вправо, оценка (36) ухудшается из-за быстрого роста функции $\omega_{n+1}(x)$. Это хорошо видно на рис. 2, где в качестве примера приведен график функции $\omega_4(x)$ с корнями $x_0 = -3/2$, $x_1 = -1/2$, $x_2 = 1/2$, $x_3 = 3/2$:

$$\omega_4(x) = \left(x^2 - \frac{9}{4} \right) \left(x^2 - \frac{1}{4} \right).$$

Ее наибольшее по модулю значение на отрезке $[-2, 2]$ равно единице. Однако уже в точках $x = \pm 2$ за пределами отрезка полином $\omega_4(x)$ принимает значение

$$g(\pm 2) = \frac{105}{16} = 6.5625.$$

Из сказанного можно сделать следующий вывод. Если $x \in [x_0, x_n]$, то множитель $|\omega_{n+1}(x)|$ не обесценивает оценку (36). Такой случай называют собственно интерполяцией $f(x)$. Противоположный случай, когда точка x лежит вне отрезка называют экстраполяцией функции $f(x)$. Отметим еще особенность новведения полинома $\omega_{n+1}(x)$ резко ухудшает оценку (36) при экстраполации. Поэтому на практике экстраполяции избегают или ограничиваются многочленами невысокой степени ($n = 1, 2$), когда рост функции $|\omega_{n+1}(x)|$ не настолько критичен.**Задача 4.**Написать мажорантную оценку для погрешности (36) при вычислении приближенного значения $\sin x$ в точке $x = \pi/4$ с помощью интерполяционного полинома второй степени $P_2(x)$ (16). Сравнить ее с погрешностью (17), подсчитанной непосредственно.

Формула для погрешности (35) принимает в данном случае вид:

$$R_2\left(\frac{\pi}{4}\right) = \frac{1}{6} (-\cos \xi) \omega_3\left(\frac{\pi}{4}\right) = \cos \xi \frac{\pi^3}{1152}, \quad 0 \leq \xi \leq \frac{\pi}{2}.$$

Она правильно определяет знак погрешности, но не позволяет вычислить ее величину, поскольку значение аргумента ξ неизвестно. Чтобы получить мажорантную оценку погрешности (36), нужно заменить $\cos \xi$ на его наибольшее значение – единицу. В результате будем иметь:

$$R_2\left(\frac{\pi}{4}\right) \leq M_{n+1} \cdot |\omega_{n+1}(x)|. \quad (44)$$

где $R_2(x) = r_2(x) \omega_{n+1}(x)$,

$$\omega_{n+1}(x) = (x - x_0)^{N_0} (x - x_1)^{N_1} \dots (x - x_n)^{N_n}, \quad n+1 = N_0 + \dots + N_n$$

и рассмотреть функцию

$$g(x) = f(x) - H_2(x) - r_2(x) \omega_{n+1}(x).$$

Применяя теорему Ролля к функции $g(x)$ и ее производным с учетом кратности корней в узлах x_i и условия $g(x_i) = 0$ приходим к формуле

$$(f(x) - H_2(x)) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{n+1}(x), \quad (43)$$

которая по существу повторяет формулу (35). С ее помощью можно написать оценку типа (36):

$$|R_2(x)| \leq \frac{M_{n+1}}{(n+1)!} |\omega_{n+1}(x)|, \quad (44)$$

где ω_{n+1} – максимальное значение модуля функции $f^{(n+1)}(x)$ (37). Здесь полином $\omega_{n+1}(x)$ (42) является производным полиномом $H_n(x)$ на слух кратных корней.

Построение полинома Эрмита в общем случае при производном числе узлов и их кратности приводят к довольно громоздким выражениям и редко используется. Поэтому мы ограничимся двумя примерами, встречающимися на практике.

Пример 1.Построить интерполяционный полином Эрмита для функции $f(x)$ по известным значениям в узлах $f_i = f(x_i)$, $k = 0, 1, \dots, m$ и значению $f'(x_j) = f'_j$ в одном из узлов x_j .Степень полинома $H_n(x)$ в данном случае равна $m+1$.Будем искать $H_{n+1}(x)$ в виде

$$H_{n+1}(x) = \sum_{k=0}^m f_k \frac{(x - x_0) \dots (x - x_k) \dots (x - x_n)}{(x_k - x_0) \dots (x_k - x_{k-1}) \dots (x_k - x_{k+1}) \dots (x_k - x_n)} \left(\frac{x - x_i}{x_i - x_j} \right) + \left[f_j + \alpha_j (x - x_i) \right] \frac{(x - x_0) \dots (x - x_{j-1}) \dots (x - x_{j+1}) \dots (x - x_n)}{(x_i - x_0) \dots (x_i - x_{j-1}) \dots (x_i - x_{j+1}) \dots (x_i - x_n)}.$$

Здесь выражение, стоящее под знаком суммы, есть обычные составляющие полинома в форме Лагранжа в узлах x_k , $k \neq j$, «усиленные» дополнительными множителями $(x - x_j)/(x_i - x_j)$. Слагаемое, отвечающее кратному узлу $x = x_j$, выделено отдельно как особое. Постоянная α_j подлежит определению.Из структуры $H_{n+1}(x)$ видно, что $H_{n+1}(x_i) = f_i$, $i = 0, 1, \dots, m$. Найдем производную $H'_{n+1}(x_i)$ в узле $x = x_j$. Слагаемые, стоящие под знаком суммы, содержат множители $(x - x_j)^2$ и потому их производные обращаются в нуль при $x = x_j$. Таким образом, $P_2(x)$ синус не только в граничных точках отрезка $\left[0, \frac{\pi}{2}\right]$, но и во внутренней точке $x = \pi/6$ приводит к тому, что полином $P_2(x)$ приближает синус на отрезке $\left[0, \frac{\pi}{2}\right]$.лучше чем полином $H_2(x)$. Это хорошо видно при сравнении рис. 1 и рис. 3. Подсчет погрешностей (17) и (51) в точке $x = \pi/4$ является дополнительным тому подтверждением.**§2. Интерполяция сплайнами.**

Увеличение степени интерполяционного полинома может оказаться невыгодным из-за быстрого роста объема вычислений. К тому же даже во время он проводит к повышению точности. Во второй половине XX века с появлением компьютеров и развитием современной вычислительной математики при обработке больших таблиц получила развитие новая методика – строительство сплайнами с помощью кусочно-полиномиальной интерполяции с использованием полиномов сравнительно низких степеней. Наиболее удобными оказались полиномы третьей степени. Такие конструкции получили название кубических сплайнов.

2.1. Определение кубического сплайна.Пусть на отрезке $[a, b]$ задана функция $y = f(x)$. Рассмотрим сетку узлов

$$a = x_0 < x_1 < x_2 < \dots < x_n = b$$

и обозначим через Δx_i расстояние между смежными узлами

$$\Delta x_i = x_{i+1} - x_i, \quad i = 1, \dots, n$$

Определение:

Назовем кусочным сплайном функцию $y = f(x)$, $x \in [a, b]$ на сетке (54) функцию $S(x)$

2.2. Формулировка системы уравнений для коэффициентов кубического сплайна

Среди задач построения сплайна к отысканию коэффициентов упомянутых полиномов третьей степени на каждом из отрезков $[x_{i-1}, x_i]$. Для этого сопоставим отрезку $[x_{i-1}, x_i]$ полином $S_i(x)$, для удобства записанный в виде:

$$S_i(x) = a_i + b_i(x - x_i) + \frac{c_i}{2}(x - x_i)^2 + \frac{d_i}{6}(x - x_i)^3, \quad x \in [x_{i-1}, x_i], \quad i = 1, \dots, n. \quad (56)$$

При этом, очевидно:

$$S'_i(x) = b_i + c_i(x - x_i) + \frac{d_i}{2}(x - x_i)^2, \quad (57)$$

$$S''_i(x) = c_i + d_i(x - x_i), \quad (58)$$

так, что

$$S_i(x_i) = a_i, \quad S'_i(x_i) = b_i, \quad S''_i(x_i) = c_i. \quad (59)$$

Для выполнения требований (S3) в узлах интерполяции с номерами $i = 1, \dots, n$ следует положить:

$$a_i = f_i(x_i), \quad i = 1, \dots, n \quad (60)$$

Требуя непрерывности сплайна в узлах x_i ($i = 1, \dots, n-1$) и выполнения условия (S3) при $i = 0$, получим:

$$S_i(x_{i-1}) = f_{i-1}, \quad i = 1, \dots, n \quad (61)$$

или

$$f_i + b_i(x_{i-1} - x_i) + \frac{c_i}{2}(x_{i-1} - x_i)^2 + \frac{d_i}{6}(x_{i-1} - x_i)^3 = f_{i-1}, \quad i = 1, \dots, n. \quad (62)$$

Это равенство можно переписать следующим образом:

$$b_i h_i - \frac{c_i}{2} h_i^2 + \frac{d_i}{6} h_i^3 = f_{i-1}, \quad i = 1, \dots, n. \quad (62)$$

Условие (S2) непрерывности первой производной $S'(x)$ в узлах x_i ($i = 1, \dots, n-1$) принимает вид:

$$S'_i(x_{i-1}) = S'_{i-1}(x_{i-1}) = b_{i-1}, \quad i = 2, \dots, n \quad (63)$$

и приводит к соотношению

$$b_i - c_i h_i + \frac{d_i}{2} h_i^2 = b_{i-1}, \quad i = 2, \dots, n \quad (64)$$

или

$$c_i h_i - \frac{d_i}{2} h_i^2 = b_i - b_{i-1}, \quad i = 2, \dots, n \quad (64)$$

Аналогичным образом условия непрерывности второй производной $S''(x)$ в тех же узлах:

$$S''_i(x_{i-1}) = S''_{i-1}(x_{i-1}) = c_{i-1}, \quad i = 2, \dots, n \quad (65)$$

означают, что

$$d_i h_i = c_i - c_{i-1}, \quad i = 2, \dots, n. \quad (66)$$

$$a_1 = 1, \quad a_2 = 3; \quad d_1 = 2, \quad d_2 = -2; \quad b_1 = 4/3, \quad b_2 = 7/3.$$

Теперь можно выписать кубические полиномы, определяющие сплайн:

$$S(x) = \begin{cases} S_1(x) = 1 + \frac{4}{3}x + x^2 + \frac{1}{3}x^3, & -1 \leq x \leq 0, \\ S_2(x) = 3 + \frac{7}{3}(x-1) - \frac{1}{3}(x-1)^3, & 0 \leq x \leq 1. \end{cases} \quad (82)$$

Легко проверить, что построенная таким образом функция $S(x)$ непрерывна вместе с первой и второй производной во внутренней узловой точке $x=0$.

В заключение вычислим значение сплайна в точке $x=1/2$, т. е. подсчитаем приближение $\sqrt{3}$:

$$\sqrt{3} \approx S_2\left(\frac{1}{2}\right) = \frac{15}{8}, \quad x = \sqrt{3} - \frac{15}{8} = -0,142949. \quad (83)$$

Значительная погрешность обусловлена прежде всего большим шагом $h=1$. Определенную роль играют также условия S4:

$$S(-1) = S'(1) = 0. \quad (84)$$

Вторая производная рассматриваемой функции $f(x) = 3^x$ в точках $x = \pm 1$ в итоге не обращается, т. е. условие (84) дает о ней искаченную информацию. Если учесть при построении сплайна истинные значения функции $f''(x)$ в точках ± 1 , то точность аппроксимации улучшится.

2.5. Сходимость и точность интерполяирования сплайнами.

При обсуждении эффективности численного метода в первую очередь обращают внимание на две характеристики:

1. Условие сходимости метода (сходимость).

Речь идет о минимальных по возможности ограничениях, при которых приближенное решение задачи стремится к точному решению задачи.

Сходимость означает, что данный метод в принципе позволяет найти решение задачи с любой степенью точности.

2. Скорость сходимости (точности).

Это характеристика близости приближенного решения к точному (характеристика скорости сходимости) при некоторых дополнительных ограничениях.

Посмотрим как решаются эти вопросы в теории сплайнов.

Итак, на сегменте $[a, b]$ задана функция $f(x)$ и построена сетка

$$a = x_0 < x_1 < x_2 < \dots < x_n = b, \quad h_i = x_i - x_{i-1} > 0.$$

Введем в рассмотрение величину

$$h = \max_{1 \leq i \leq n} h_i. \quad (85)$$

Приведем без доказательства две теоремы.

Теорема 1. Пусть $f(x)$ непрерывна на сегменте $[a, b]$, тогда для любого $\varepsilon > 0$ можно указать $\delta(\varepsilon) > 0$ такое, что при любой сетке, удовлетворяющей условию $h < \delta$ справедливо неравенство

$$|f(x) - S(x)| < \varepsilon \quad \forall x \in [a, b]. \quad (86)$$

Квадрат погрешности в точке $x = x_i$ для функции $F(x)$ (91) с коэффициентами (100) можно записать в виде

$$\begin{aligned} \delta_i^2 &= \left\{ y_i - \sum_{k=0}^n (\bar{a}_k + \Delta a_k) \phi_k(x_i) \right\}^2 = \\ &= \left\{ \left[y_i - \sum_{k=0}^n \bar{a}_k \phi_k(x_i) \right] - \sum_{k=0}^n \Delta a_k \phi_k(x_i) \right\}^2 = \left\{ y_i - \sum_{k=0}^n \bar{a}_k \phi_k(x_i) \right\}^2 - \\ &\quad - 2 \left\{ y_i - \sum_{k=0}^n \bar{a}_k \phi_k(x_i) \right\} \left(\sum_{k=0}^n \Delta a_k \phi_k(x_i) \right) + \left(\sum_{k=0}^n \Delta a_k \phi_k(x_i) \right)^2. \end{aligned} \quad (101)$$

Здесь в среднем слагаемым мы заменили в одной из сумм индекс суммирования k на l , чтобы не использовать один и тот же индекс в двух разных суммах и иметь возможность перемножить их поочередно.

Чтобы получить суммарную квадратичную погрешность, нужно просуммировать выражение (101) для δ_i^2 по индексу i . Первые слагаемые не содержат Δa_k . Их сумма дает погрешность J , вычисленную для функции (91) с коэффициентами (99) \bar{a}_k .

Рассмотрим теперь сумму вторых слагаемых, которые зависят от Δa_k линейно:

$$\begin{aligned} -2 \sum_{i=0}^n &\left\{ \left[y_i - \sum_{k=0}^n \bar{a}_k \phi_k(x_i) \right] - \sum_{k=0}^n \Delta a_k \phi_k(x_i) \right\} = \\ &= -2 \sum_{i=0}^n \Delta a_i \left\{ \sum_{k=0}^n y_k \phi_k(x_i) - \sum_{k=0}^n \bar{a}_k \phi_k(x_i) \phi_k(x_i) \right\} = \\ &= -2 \sum_{i=0}^n \Delta a_i \left\{ b_i - \sum_{k=0}^n \bar{a}_k y_k \right\} = 0. \end{aligned} \quad (102)$$

Здесь мы поменяли местами порядок суммирования и воспользовались тем, что коэффициенты \bar{a}_k удовлетворяют системе уравнений (95).

С учетом (102) будем иметь

$$\begin{aligned} &(\bar{a}_0 + \Delta a_0, \bar{a}_1 + \Delta a_1, \dots, \bar{a}_n + \Delta a_n) = \\ &= J(\bar{a}_0, \bar{a}_1, \dots, \bar{a}_n) + \sum_{i=0}^n \sum_{k=0}^n \Delta a_i \phi_k(x_i) > J(\bar{a}_0, \bar{a}_1, \dots, \bar{a}_n). \end{aligned} \quad (103)$$

Формула (103) показывает, что функция $F(x)$ (91) с коэффициентами \bar{a}_k (100), полученным в результате решения уравнений (95), действительно минимизирует суммарную квадратичную погрешность J . Если мы возьмем любой другой набор коэффициентов (100), отличный от (99), то согласно формуле (103) к погрешности $u(\bar{a}_0, \bar{a}_1, \dots, \bar{a}_n)$ добавится положительное слагаемое и она увеличится.

Итак, чтобы построить наилучшее приближение (91) сеточной функции (1), (2) по методу наименьших квадратов, нужно взять в качестве коэффициентов разложения a_k решение системы линейных уравнений (95).

Наконец, дополнительные граничные условия (S4) дают еще два уравнения

$$\begin{cases} S'_1(x_0) = S'_1(a) = c_1 - d_1 h_1 = 0 \\ S''_n(x_n) = S''_n(b) = c_n = 0 \end{cases}. \quad (67)$$

В итоге мы получили замкнутую систему (62), (64), (66), (67), содержащую в сумме 3n линейных уравнений для отыскания 3n неизвестных: $b_i, c_i, d_i, i = 1, 2, \dots, n$

2.3. Редукция системы.

Удобно формально ввести еще одно неизвестное c_0 , положив при этом $c_0 = 0$, и первое уравнение в (67) переписать в виде:

$$d_1 h_1 = c_1 - c_0,$$

то есть в форме аналогичной (66).

Теперь уравнения (66) и (67) естественно представить в единобразном виде

$$d_i h_i = c_i - c_{i-1}, \quad i = 1, 2, \dots, n$$

$$c_0 = 0, \quad c_n = 0. \quad (69)$$

Обратим внимание на то, что из системы (68) можно выразить все коэффициенты d_i через разности $c_i - c_{i-1}$, а затем из системы (62) выразить через c_i и c_{i-1} коэффициенты b_i . Подставляя полученные выражения в (64), придем к системе линейных уравнений для c_i :

$$\frac{1}{3} c_{i-2} h_{i-1} + \frac{2}{3} c_{i-1} (h_{i-1} + h_i) + \frac{1}{3} c_i h_i = 2 \left(\frac{f_{i-1} - f_i}{h_{i-1}} - \frac{f_i - f_{i+1}}{h_i} \right), \quad i = 2, 3, \dots, n. \quad (70)$$

Сдвигая индекс i на единицу, получим симметричную форму записи уравнений (70):

$$h_{i-1} c_{i-1} + 2(h_{i-1} + h_i) c_i + h_i c_{i+1} = 6 \left(\frac{f_{i-1} - f_i}{h_{i-1}} - \frac{f_i - f_{i+1}}{h_i} \right), \quad i = 1, \dots, n-1. \quad (71)$$

Кроме того, согласно (69)

$$c_0 = c_n = 0. \quad (72)$$

Система (71) содержит $n-1$ уравнение с $(n-1)$ -ой неизвестной: c_1, c_2, \dots, c_{n-1} . Величины c_i и c_{i-1} определены дополнительными соотношениями (72). Если сетка (54) равномерная, т. е. $h_i = h = \text{const}$, то уравнения (71) принимают особенно простой вид:

$$c_{i-1} + 4c_i + c_{i+1} = \frac{f_{i-1} - 2f_i + f_{i+1}}{h^2}. \quad (73)$$

Для уравнений системы (71) выполняется условие диагонального преобразования. Отсюда следует существование и единственность решения задачи (71), (72). По найденным величинам c_i можно рассчитать остальные коэффициенты сплайна по формулам

$$d_i = \frac{c_{i-1} - c_{i+1}}{h_i}, \quad i = 1, \dots, n \quad (74)$$

и

$$b_i = \frac{1}{2} h_i c_i + \frac{h}{6} h_i^2 d_i + \frac{f_i - f_{i+1}}{h}, \quad i = 1, \dots, n, \quad (75)$$

завершив тем самым построение сплайна. Теорема доказана.

2.4. Замечание о решении системы.

Уравнения (71) имеют так называемую трехточечную структуру, общий вид таких систем

$$A_j y_{i-1} + C_j y_i + B_j y_{i+1} = f_i, \quad i = 1, 2, \dots, n-1, \quad (76)$$

$$y_0 = 0, \quad y_n = 0. \quad (77)$$

соответствует системе линейных уравнений с треугольной матрицей T для определения вектора неизвестных $y = (y_1, y_2, \dots, y_{n-1})$:

$$Ty = F, \quad (78)$$

где

$$T = \begin{vmatrix} C_1 & B_1 & 0 & 0 & 0 & 0 \\ A_2 & C_2 & B_2 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & A_n & C_n & B_n & 0 & 0 \end{vmatrix}, \quad F = \begin{vmatrix} F_1 \\ F_2 \\ \vdots \\ F_{n-1} \\ F_n \end{vmatrix}. \quad (79)$$

При этом легко видеть, что в нашем случае

$$|C_i| > |A_i| + |B_i|, \quad i = 1, \dots, n-1, \quad (79)$$

поскольку

$$C_i = 2(h_i + h_{i+1}), \quad A_i = h_i, \quad B_i = h_{i+1}. \quad (80)$$

Как было показано в главе 1, решение подобных систем эффективно осуществляется методом прогонки.

Задача 6.

Рассмотреть функцию $y = f(x) = 3^x$ на отрезке $[-1, 1]$ с узами интерполяции $x_0 = -1$, $x_1 = 0$, $x_2 = 1$. Построить кубический сплайн. Найти его значение при $x = 1/2$, т. е. вычислить приближение $\sqrt{3}$. Подсчитать погрешность.

В рассматриваемомслучае мы имеем равномерную сетку с шагом $h = 1$. У нас одна внутренняя точка x_1 и две граничные — x_0 и x_2 . Система (73) сводится к одному уравнению относительно коэффициента c_1 , которое с учетом дополнительных соотношений (70), определяющих нулевые значения коэффициентов c_0 и c_2 , принимает вид:

$$4c_1 = 6 \left(\frac{1}{3} - 2 + 3 \right). \quad (81)$$

Таким образом, в нашей задаче:

$$c_0 = 0, \quad c_2 = 2, \quad c_1 = 0.$$

Остальные коэффициенты сплайна находятся по формулам (60), (74), (75):

иными словами $S_h(x)$ при $h \rightarrow 0$ равномерно сходится к непрерывной функции $f(x)$.

Теорема 2. Пусть $f(x)$ имеет на сегменте $[a, b]$ четыре непрерывные производные и дополнительно удовлетворяет условию $f''''(a) = f''''(b) = 0$. Тогда имеют место неравенства (оценки):

$$|f(x) - S(x)| \leq M_1 h^4 \quad \forall x \in [a, b], \quad (87)$$

$$|f'(x) - S'(x)| \leq M_2 h^3 \quad \forall x \in [a, b], \quad (88)$$

$$|f''(x) - S''(x)| \leq M_3 h^2 \quad \forall x \in [a, b], \quad (89)$$

$$M_4 = \max_{[a,b]} |f^{(4)}(x)|. \quad (90)$$

Функция $F(x)$ (91) с набором коэффициентов, удовлетворяющих этому требованию, называют наилучшим приближением по методу наименьших квадратов.

Построение наилучшего приближения сводится к классической задаче математического анализа об экстремуме функции нескольких переменных. Метод решения этой задачи известен. Необходимым условием экстремума является равенство нулю в экстремальном точке всех первых частных производных рассматриваемой функции. В случае (93) это дает

$$\frac{\partial J}{\partial a_i} = 2 \sum_{j=0}^n (y_j - \sum_{k=0}^n a_k \phi_k(x_j)) \phi_i(x_j) = 0 \quad i = 0, 1, \dots, m. \quad (94)$$

Оставим члены, содержащие a_m , слева и поменяем в них порядок суммирования по индексам i и k . Члены, содержащие y_i , перенесем направо. В результате уравнения (94) примут вид

$$\sum_{k=0}^m y_k a_k = b_i, \quad i = 0, 1, \dots, m, \quad (95)$$

где

$$y_k = \sum_{i=0}^n \phi_i(x_j) \phi_k(x_i), \quad (96)$$

$$b_i = \sum_{j=0}^n \phi_i(x_j) y_j. \quad (97)$$

Мы получили систему линейных алгебраических уравнений (95), в которой роль неизвестных играют искомые коэффициенты разложения a_0, a_1, \dots, a_m . Число уравнений и число неизвестных в этой системе совпадают и равно $m+1$. Матрица коэффициентов системы G состоит из элементов g_{ij} , которые определяются формулой (96). Ее называют матрицей Грама для системы функций $\phi_0(x), \phi_1(x), \dots, \phi_m(x)$ на сетке (1). Отметим, что матрица Грама является симметричной: для ее элементов, согласно (96), справедливо равенство $g_{ij} = g_{ji}$. Числа b_i , стоящие в правой части уравнений (95), вычисляются по формуле (97) через значения y_i сеточной функции (2).

Предположим, что функции $\phi_0(x), \phi_1(x), \dots, \phi_m(x)$ выбраны такими, что определитель матрицы Грама, отличен от нуля:

$$\Delta = \det G \neq 0. \quad (98)$$

В этом случае при любой правой части система (95) имеет единственное решение:

$$\bar{a}_0, \bar{a}_1, \dots, \bar{a}_m. \quad (99)$$

Глава 3. ЧИСЛЕННОЕ ИНТЕГРИРОВАНИЕ

§1. Формула Ньютона-Лейбница и численное интегрирование.
Из курса математического анализа Вы знакомы с вычислением определенных интегралов с помощью формулы Ньютона-Лейбница:

$$I = \int_a^b f(x) dx = F(b) - F(a). \quad (1)$$

где $F(x)$ – любая первообразная подынтегральной функции $f(x)$ на отрезке $[a, b]$. Формула Ньютона-Лейбница играет важную роль, устанавливая связь задачи определенного интегрирования с задачей отыскания первообразной (с задачей неопределенного интегрирования). Она позволяет вычислять интегралы от функций, первообразные которых тоже являются элементарными функциями. Например,

$$\int \frac{dx}{x} = \ln|x| \Big|_a^b = \ln b - \ln a. \quad (2)$$

Однако существует много простых функций, первообразные которых не выражаются через элементарные функции. В качестве примера можно привести такие функции как e^{-x^2} или $\sin x^2$. Для них описанный способ вычисления определенных интегралов не применим. Формула Ньютона-Лейбница не позволяет также вычислять интегралы от функций, которые задаются графиком или таблицей. Иными словами, она не дает общего, универсального метода нахождения определенного интеграла от произвольной функции $f(x)$ по ее значениям на отрезке $[a, b]$, она не является алгоритмом решения рассматриваемой задачи.

Универсальные алгоритмы вычисления определенных интегралов дают формулы численного интегрирования или, как их обычно называют, квадратурные формулы (буквально формулы вычисления площадей). Квадратурные формулы имеют вид:

$$I = \int_a^b f(x) dx = \sum_{i=1}^n c_i f(x_i) + R_n. \quad (3)$$

Здесь точки $x_i \in [a, b]$ называются узлами, коэффициенты c_i – весовыми множителями или просто весами, величина R_n – остаточным членом или погрешностью. Узлы и веса подбираются таким образом, чтобы выполнялось предельное равенство:

$$\lim_{n \rightarrow \infty} R_n = 0, \text{ так что } \lim_{n \rightarrow \infty} \sum_{i=1}^n c_i f(x_i) = I. \quad (4)$$

Суть этого требования заключается в следующем. Если пренебречь в формуле (3) остаточным членом R_n , то получится приближенное равенство:

$$I = \int_a^b f(x) dx \approx \sum_{i=1}^n c_i f(x_i). \quad (5)$$

$$I = \int_a^b f(x) dx = T_n + \beta_n, \quad (14)$$

В квадратурной формуле (14) узлами являются точки x_i (6). Все весовые коэффициенты, кроме x_0 и x_n , равны $h = (b-a)/n$, а весовые коэффициенты при $i=0$ и $i=n$ имеют значения в два раза меньше. Для остаточного члена введено специальное обозначение β_n . Формулу (14) называют квадратурной формулой трапеций. С точностью до β_n она выражает площадь криволинейной трапеции, соответствующую интегралу I , через сумму площадей трапеций (12) (см. рис. 2).

Формула (8) для величины P_n изначально строилась как интегральная сумма. При выводе формулы (13) для величины T_n понятие интегральной суммы не использовалось. Однако теперь, когда формула уже получена, видно, что величину T_n тоже можно интерпретировать как интегральную сумму. Чтобы убедиться в этом, рассмотрим разбиение отрезка $[a, b]$ на частичные отрезки точками ξ_i (7). Оно дает $n+1$ отрезок. Две крайние $[\xi_0, \xi_1]$ и $[\xi_n, \xi]$ имеют длину $h/2$, а остальные – длину h . Выберем для образования интегральной суммы на крайних отрезках значения функции $f(x)$ в точках a и b , а на остальных отрезках $[\xi_i, \xi_{i+1}]$ – значения функции $f(x)$ в их средних точках x_i ($1 \leq i \leq n-1$). Образованная таким образом интегральная сумма соответствует выражению (13) для T_n .

Выход квадратурной формулы Симпсона развивается описанной подольше. Теперь для аппроксимации функции $f(x)$ используется не кусочно – линейное, а кусочно – квадратичное интерполирование.

Будем считать n четным и группируем отрезки $[x_{i-1}, x_i]$ парами: первая пара $[a, x_1], [x_1, x_2]$, вторая пара $[x_2, x_3], [x_3, x_4]$ и т. д. Для каждого двойного отрезка $[x_{2j-2}, x_{2j}]$ построим интерполяционный полином второй степени в форме Лагранжа, принимающий в узлах $x_{2j-2}, x_{2j-1}, x_{2j}$ значения функции $f(x)$. В результате получим аппроксимирующую функцию $g_n(x)$ на отрезке $[a, b]$ в виде кусочно – квадратичной функции:

$$g_n(x) = f(x_{2j-1}) \frac{(x-x_{2j-1})(x-x_{2j})}{2h^2} + f(x_{2j}) \frac{(x-x_{2j-1})(x-x_{2j})}{(-h)^2} + f(x_{2j}) \frac{(x-x_{2j-1})(x-x_{2j})}{2h^2}, \quad x \in [x_{2j-2}, x_{2j}], \quad 1 \leq j \leq n/2. \quad (15)$$

Пронтегрировав полином второй степени (15) по отрезку $[x_{2j-2}, x_{2j}]$, получим

где η_j^* и η_j^{**} – некоторые точки отрезка $[x_{i-1}, x_i]$. Существование таких точек гарантировано, но их точное положение неизвестно. (См. В. А. Ильин, Э. Г. Позник «Основы математического анализа». М. 1965. С. 389–397.)

Суммируя равенства (25) и (26) по i , получим формулы (9) и (14) со следующими выражениями для остаточных членов

$$\alpha_n = \frac{h}{24} \sum_{i=1}^n f''(\eta_j^*), \quad (27)$$

$$\beta_n = -\frac{h^3}{12} \sum_{i=1}^n f''(\eta_j^{**}). \quad (28)$$

Рассмотрим суммы

$$h \sum_{i=1}^n f''(\eta_j^*) \text{ и } h \sum_{i=1}^n f''(\eta_j^{**}). \quad (29)$$

Функция $f''(x)$ по предположению непрерывна и, следовательно, интегрируема на отрезке $[a, b]$. С учетом этого замечания выражения (29) можно рассматривать как

интегральные суммы для интеграла $\int_a^b f''(x) dx$. Отсюда следует вывод:

$$\lim_{n \rightarrow \infty} h \sum_{i=1}^n f''(\eta_j^*) = \int_a^b f''(x) dx = f'(b) - f'(a), \quad (30)$$

$$\lim_{n \rightarrow \infty} h \sum_{i=1}^n f''(\eta_j^{**}) = \int_a^b f''(x) dx = f'(b) - f'(a). \quad (31)$$

Предельные равенства (30) и (31) позволяют записать остаточные члены квадратурных формул прямоугольников и трапеций в виде

$$\alpha_n = \frac{1}{n^2} (A + \mu_n), \quad (32)$$

$$\beta_n = \frac{1}{n^2} (B + \nu_n), \quad (33)$$

где

$$A = \frac{(b-a)^3}{24} (f'(b) - f'(a)), \quad (34)$$

$$\mu_n = \frac{(b-a)^3}{24} \left(\int_a^b f''(\eta_j^*) - \int_a^b f''(x) dx \right) \rightarrow 0, \text{ при } n \rightarrow \infty, \quad (35)$$

$$B = \frac{(b-a)^3}{12} (f'(b) - f'(a)), \quad (36)$$

$$\nu_n = -\frac{(b-a)^3}{12} \left(\int_a^b f''(\eta_j^{**}) - \int_a^b f''(x) dx \right) \rightarrow 0, \text{ при } n \rightarrow \infty. \quad (37)$$

Формулы (32) и (33) выделяют в остаточных членах главные слагаемые A/n^2 и B/n^2 , которые при возрастании n стремятся к нулю как n^{-2} . Важно подчеркнуть, что

Условие (4), которое называют сходимостью, позволяет сделать нигде непрерывную в равенстве (5) меньше любого наперед заданного числа за счет выбора достаточно большого n . Таким образом, открывается возможность вычислить интеграл I с любой наперед заданной точностью, по значениям функции $f(x)$, взятым в разных точках x , отрезка $[a, b]$. Чем выше требование точности, тем больше слагаемых приходится уделять в сумме. За точность приходится платить увеличением объема вычислений.

В заключение сделаем следующее замечание. Подставляя в формулу (3) функцию $f(x) = 1$, получим:

$$(b-a) = \sum_{i=1}^n c_i + R_n.$$

Обычно весовые коэффициенты c_i подбираются таким образом, чтобы выполнялось равенство:

$$(b-a) = \sum_{i=1}^n c_i,$$

т. е., чтобы при интегрировании константы равенство (5) было не приближенным, а точным.

В следующих параграфах этой главы мы обсудим методы построения квадратурных формул с разными сторонами разбивкой оценки их точности.

§2. Квадратурные формулы прямоугольников, трапеций, Симпсона.

2.1. Квадратурные формулы прямоугольников, трапеций, Симпсона и их особенности.

С квадратурными формулами прямоугольников, трапеций, Симпсона Вы уже встречались в курсе математического анализа, поэтому их вывод будет изложен кратко.

Возьмем произвольное целое число n и разбъем отрезок $[a, b]$, по которому ведется интегрирование, на n равных отрезков длиной $h = (b-a)/n$ по

точкам

$$x_i = a + ih, \quad 0 \leq i \leq n.$$

Для дальнейшего нам также понадобится среднее значение этих отрезков

$$\xi_i = a + (i-1/2)h, \quad \xi_i \in [x_{i-1}, x_i], \quad 1 \leq i \leq n. \quad (7)$$

Идея вывода формулы прямоугольников очень проста. Построим с помощью проведенного разбиения интегральную сумму, в которой значения функции $f(x)$ для каждого отрезка $[x_{i-1}, x_i]$ вычисляются в ее средней точке ξ_i (7):

$$P_n = \frac{b-a}{n} \sum_{i=1}^n f(\xi_i). \quad (8)$$

Примите внимание то, что интегральная сумма дает приближенное значение интеграла, можно написать:

$$I = P_n + \alpha_n. \quad (9)$$

$$\int_{x_{j-1}}^{x_j} g_n(x) dx = \frac{h}{3} \{ f(x_{j-1}) + 4f(x_{j-1}) + f(x_j) \}, \quad h = \frac{b-a}{n}. \quad (16)$$

Интеграл от функции $g_n(x)$ по всему отрезку $[a, b]$ равен сумме интегралов (16)

$$S_n = \int_a^b g_n(x) dx = \sum_{j=1}^n \int_{x_{j-1}}^{x_j} g_n(x) dx =$$

$$= \frac{b-a}{3n} \{ f(a) + 4f(x_1) + 2f(x_2) + \dots + 2f(x_{n-2}) + 4f(x_{n-1}) + f(b) \}. \quad (17)$$

(Напомним, что число n должно быть обязательно четным). Величина S_n (17) дает приближенное значение интеграла I :

$$I = \int_a^b f(x) dx = S_n + \gamma_n. \quad (18)$$

Узлами квадратурной формулы (17), как и формулы трапеций (14), являются точки x_i (6). Весовые коэффициенты в узлах с четными и нечетными номерами имеют разные значения. Для остаточного члена введено обозначение γ_n . Формула (18) называется квадратурной формулой Симпсона.

Представление (17) для S_n как и представление (13) для T_n , также можно рассматривать как интегральную сумму. Для ее построения нужно разбить отрезок $[a, b]$ на $(n+1)$ частичный отрезок с помощью n внутренних точек

$$\eta_{j-1} = x_{j-1} - 2h/3, \quad \eta_j = x_{j-1} + 2h/3, \quad 1 \leq j \leq n/2 \quad (19)$$

и двух граничных точек

$$\eta_0 = a \text{ и } \eta_{n+1} = b. \quad (20)$$

В результате получаются отрезки $[\eta_0, \eta_1], [\eta_1, \eta_2], \dots, [\eta_n, \eta_{n+1}]$ и $[\eta_{n+1}, b]$ имеют длину $h/3$. Отрезки, в центре которых лежат точки x_i с четными номерами, – длину $2h/3$, отрезки, в центре которых лежат точки x_i с нечетными номерами, – длину $4h/3$.

Для построения интегральной суммы, соответствующей данному разбиению, возьмем для крайних отрезков значения функции $f(x)$ в точках a и b , для остальных отрезков – значение функции $f(x)$ в их средних точках x_i . В результате получим интегральную сумму в виде выражения (17). Разные длины частичных отрезков приводят к своеобразному чередованию коэффициентов в виде двоек, четверок и единиц в крайних точках.

Заканчивая обсуждение формул (13) для T_n и (17) для S_n , установим полезную для дальнейшего связь между этими величинами

$$S_n = \frac{4}{3} T_n - \frac{1}{3} \gamma_n. \quad (21)$$

Здесь T_n – сумма (13) с вдвое меньшим числом слагаемых и, соответственно, сдвое большим шагом. Благодаря этому при ее образовании в качестве узлов используются точки x_i (6) только с четными номерами. Поскольку в формуле Симпсона n

коэффициенты A (34) и B (36) от n не зависят. Дополнительные слагаемые μ_n/n^2 и ν_n/n^2 являются бесконечно малыми более высокого порядка. Если ими пренебречь по сравнению с главными слагаемыми, то получатся простые асимптотические представления остаточных членов:

$$\alpha_n \approx An^2 \text{ и } \beta_n \approx Bn^2. \quad (38)$$

Их относительная точность возрастает при увеличении n .

Теперь получим другое представление остаточных членов. Из курса математического анализа известно следующее утверждение:

Лемма. Пусть функция $\phi(x)$ непрерывна на отрезке $[a, b]$ и пусть x_1, x_2, \dots, x_n – некоторые точки этого отрезка. Тогда на отрезке $[a, b]$ найдется такая точка η , что

$$\sum_{i=1}^n \phi(x_i) = \phi(\eta), |\beta_i| \leq \frac{(b-a)^2 M_i}{12n^2}. \quad (39)$$

Иными словами, среднее арифметическое значение непрерывной функции в нескольких точках отрезка $[a, b]$, равно ее значению в одной из точек этого отрезка.

Приимется это утверждение к суммам (27) и (28), получим другое представление остаточных членов α_n и β_n :

$$\alpha_n = \frac{(b-a)^3}{24n^2} f''(\eta^*), \quad \eta^* \leq \eta \leq b, \quad (40)$$

$$\beta_n = \frac{(b-a)^3}{12n^2} f''(\eta^{**}), \quad \eta^{**} \leq \eta \leq b. \quad (41)$$

Формулы (40) и (41) не позволяют вычислить остаточные члены: существование точек η^* и η^{**} на отрезке $[a, b]$ гарантировано, но их положение неизвестно. Однако эти формулы можно использовать для оценки остаточных членов. Пусть известно число M_2 , которое является максимумом для второй производной функции $f(x)$:

$$|f''(x)| \leq M_2, \quad a \leq x \leq b, \quad (42)$$

тогда равенства (40) и (41) можно заменить неравенствами:

$$|\alpha_n| \leq \frac{(b-a)^3 M_2}{24n^2}, \quad (43)$$

$$|\beta_n| \leq \frac{(b-a)^3 M_2}{12n^2}. \quad (44)$$

При заданной точности ε они позволяют определить число узлов n , которое нужно использовать при вычислении интеграла по рассматриваемым квадратурным формулам.

В случае, когда вторая производная функции $f(x)$ является знаконаправленной на отрезке $[a, b]$, формулы (40) и (41) позволяют определить знаки остаточных членов. При этом существенно то, что они оказываются противоположными. Пусть, например,

В квадратурной формуле (9) узлами являются точки ξ_i (7), все весовые множители одинаковы и равны $h = (b-a)/n$. Для остаточного члена введен специальный обозначение α_n .

Формулу (9) называют формулой прямоугольников. Причина такого названия имеет простой геометрический смысл. Величина P_n (8) представляет собой сумму площадей трапеций, соответствующих исходному интегралу (см. рис. 1).

Идея вывода квадратурных формул трапеций и Симпсона иная. Она заключается в том, чтобы сопоставить подынтегральную функцию $f(x)$ близкую ей функцию $g_n(x)$, которую можно проинтегрировать, и приближенно заменить искомый интеграл I интегрированием от этой функции.

Рассмотрим, как данная идея реализуется при выводе формулы трапеций. В этом случае в качестве аппроксимирующей функции $g_n(x)$ берется кусочно – линейная функция. На каждом из частичных сегментов $[x_{i-1}, x_i]$ задается формулой

$$g_n(x) = f(x_{i-1}) + \frac{f(x_i) - f(x_{i-1})}{h} (x - x_{i-1}), \quad x \in [x_{i-1}, x_i], \quad 1 \leq i \leq n. \quad (10)$$

В граничных точках отрезка $x = x_{i-1}$ и $x = x_i$ функция $g_n(x)$ принимает те же значения, что и функция $f(x)$:

$$g_n(x_{i-1}) = f(x_{i-1}), \quad g_n(x_i) = f(x_i), \quad (11)$$

т. е. она осуществляет кусочно – линейную аппроксимацию функции $f(x)$ на отрезке $[a, b]$ (см. рис. 2).

Вычислим интеграл:

$$\int_{x_{i-1}}^{x_i} g_n(x) dx = \int_{x_{i-1}}^{x_i} \left(f(x_{i-1}) + \frac{f(x_i) - f(x_{i-1})}{h} (x - x_{i-1}) \right) dx = \frac{h}{2} (f(x_{i-1}) + f(x_i)). \quad (12)$$

Этот результат имеет простой геометрический смысл: фигура ограниченная снизу отрезком $[x_{i-1}, x_i]$ оси x , сверху отрезком прямой (10), с боков вертикальными прямыми $x = x_{i-1}$ и $x = x_i$, представляет собой трапецию, площадь которой является формулой (12).

Интеграл от функции $g_n(x)$ является искомым интегралом I :

$$T_n = \sum_{i=1}^n g_n(x) dx = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} g_n(x) dx =$$

$$= \frac{b-a}{2} \left\{ \frac{1}{2} f(a) + f(x_1) + f(x_2) + \dots + f(x_{n-1}) + \frac{1}{2} f(b) \right\}. \quad (13)$$

Он дает приближенное значение интеграла I :

предполагается обязательно четным, то $n/2$ – целое число, так что выражение T_n определено.

Соотношение (21) проверяется «в лоб». Из (13) следует, что:

$$\frac{4}{3} T_n = \frac{b-a}{3} (2f(a) + 4f(x_1) + 4f(x_2) + \dots + 4f(x_{n-2}) + 4f(x_{n-1}) + 2f(b)),$$

<math display

$$|\zeta_n| \leq \frac{1}{96} \left(\frac{\pi}{2} \right)^3 < 0.041, \quad |\beta_n| \leq \frac{1}{48} \left(\frac{\pi}{2} \right)^3 < 0.081. \quad (51)$$

Перейдем к обсуждению остаточного члена γ_n в методе Симпсона, которое проводим при предположении о четырехкратной непрерывной дифференцируемости подинтегральной функции $f(x)$. Напомним, что в методе Симпсона число точек n выбирается четным, так что $n/2$ является целым числом.

Рассмотрим отрезок двойной длины $[2x]$, расположенный между точками разбиения (6) с четными номерами $[x_{j-2}, x_j]$, $1 \leq j \leq n/2$. В курсе математического анализа выводится формула:

$$\int_{x_{j-2}}^{x_j} f(x) dx = \frac{h}{3} \left[f(x_{j-2}) + 4f(x_{j-1}) + f(x_j) \right] - \frac{h^3}{90} f^{(4)}(\eta_j), \quad (52)$$

где $\eta_j \in [x_{j-2}, x_j]$. Существование такой точки гарантировано, но ее точное положение на отрезке неизвестно.

Суммируя равенства (52) по j , получим квадратурную формулу (18) со следующим выражением для остаточного члена:

$$\gamma_n = \frac{h^3}{90} \sum_{j=1}^{n/2} f^{(4)}(\eta_j). \quad (53)$$

Из формулы (53), аналогичной формулам (27), (28), можно вывести различные представления остаточного члена и изучить его свойства.

Рассмотрим сумму

$$2h \sum_{j=1}^{n/2} f^{(4)}(\eta_j). \quad (54)$$

Функция $f^{(4)}(x)$ предполагается непрерывной и, следовательно, интегрируемой на отрезке $[a, b]$. С учетом этого сумму (54) можно рассматривать как интегральную сумму для интеграла $\int_a^b f^{(4)}(x) dx$. Отсюда следует вывод

$$\lim_{n \rightarrow \infty} 2h \sum_{j=1}^{n/2} f^{(4)}(\eta_j) = \int_a^b f^{(4)}(x) dx = f''(b) - f''(a). \quad (55)$$

Предельное равенство (55) позволяет записать остаточный член квадратурной формулы Симпсона (53) в виде

$$\gamma_n = \frac{1}{n!} (C + \sigma_n), \quad (56)$$

$$C = \frac{(b-a)^4}{180} \{f''(b) - f''(a)\}, \quad (57)$$

$$\sigma_n = \frac{(b-a)^4}{180} \left\{ 2h \sum_{j=1}^{n/2} f^{(4)}(\eta_j) - \int_a^b f^{(4)}(x) dx \right\} \rightarrow 0, \text{ при } n \rightarrow \infty. \quad (58)$$

$$\alpha_n = \frac{1}{n!} (A + \mu_n) = \frac{1}{3} (P_n - P_{n/2}) + \frac{4}{3n!} (A_n - \mu_{n/2}). \quad (59)$$

Первый член в правой части этого представления остаточного члена нам известен из результатов вычислений. Он является главным. Второй член неизвестен, но он, по сравнению с первым, представляет собой бесконечно малую более высокого порядка. Если им пренебречь, то для погрешности получится простая асимптотическая формула:

$$\alpha_n \approx \frac{1}{3} (P_n - P_{n/2}). \quad (60)$$

Ее относительная точность возрастает при увеличении n .

Аналогичные формулы имеют место для погрешности метода трапеций

$$\beta_n = \frac{1}{3} (T_n - T_{n/2}) + \frac{4}{3n!} (V_n - V_{n/2}) \approx \frac{1}{3} (T_n - T_{n/2}). \quad (61)$$

Для метода Симпсона, который является методом четвертого порядка, формулы немного изменяются. Теперь соотношения, аналогичные (64), будут иметь вид:

$$I = S_{n/2} + \frac{16}{n!} (C + \sigma_n), \quad (62)$$

$$I = S_n + \frac{1}{n!} (C + \sigma_n). \quad (63)$$

(Здесь число n предполагается кратным четырем, так что $n/2$ четное число.) Проводя в (68) вычитание второй строки из первой, получим

$$\gamma_n = \frac{1}{n!} (C + \sigma_n) = \frac{1}{15} (S_n - S_{n/2}) + \frac{16}{15n!} (\sigma_n - \sigma_{n/2}). \quad (64)$$

Здесь опять член в правой части равенства известен из вычислений. Он является главным. Второй член неизвестен, но он представляет собой бесконечно малую более высокого порядка по сравнению с первым. Если им пренебречь, то получим асимптотическую формулу для приближенного вычисления погрешности по результатам двух вычислений

$$\gamma_n \approx \frac{1}{15} (S_n - S_{n/2}). \quad (65)$$

Ее относительная точность возрастает с увеличением n .

Очевидно апостериорные оценки погрешности с помощью асимптотических формул (66), (67), (70) включают в компьютерные программы численного интегрирования. Они служат критериям для завершения вычислений после того, как погрешность достигнута.

В заключение отметим следующее. Можем подставить полученные выражения для остаточных членов (65), (67), (70) в исходные квадратурные формулы (9), (14) и (18). В результате они примут вид:

$$I = \frac{4}{3} P_n - \frac{1}{3} P_{n/2} + \tilde{\alpha}_n, \quad (71)$$

$$I = \frac{4}{3} T_n - \frac{1}{3} T_{n/2} + \tilde{\beta}_n, \quad (72)$$

Таким образом, сумма весовых коэффициентов в квадратурной формуле Гаусса при любом n равна двум.

Задача 5.

Составить и решить систему уравнений (81) для квадратурной формулы Гаусса с одним узлом.

В этом случае в задаче подлежат определению два параметра: узел x_1 и весовой коэффициент c_1 . Система уравнений для их определения получается из (81) при $m=0$ и $m=1$:

$$\begin{cases} c_1 = 2 \\ c_1 x_1 = 0 \end{cases}$$

Ее решение имеет вид: $x_1 = 0$, $c_1 = 2$, так что искомая квадратурная формула записывается следующим образом:

$$\int_{-1}^1 f(x) dx = 2f(0) + \delta_1. \quad (83)$$

Выбор в качестве единственного узла средней точки отрезка $[-1, 1]$ выглядит по соображениям симметрии вполне естественно. Требование, чтобы сумма весовых коэффициентов равнялась двум (D), определяет в данном случае единственный весовой коэффициент c_1 . Квадратурная формула (83) является точной для любой линейной функции $Q_1 = a_0 + a_1 x$.

3.2. Полиномы Лежандра

Мы решим систему уравнений (81) при $n=1$. Однако решить ее «в лоб» в общем случае при произвольном n сложно. Поэтому мы будем вынуждены воспользоваться обратным путем. Для этой цели нам понадобятся полиномы Лежандра, с которыми Вы уже встречались в курсе линейной алгебры. Они определяются формулами

$$P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n. \quad (84)$$

Выпишем, используя эту формулу, несколько первых полиномов Лежандра

$$P_0(x) = 1, \quad P_1(x) = x, \quad P_2(x) = \frac{3}{2}x^2 - \frac{1}{2}, \quad P_3(x) = \frac{5}{2}x^3 - \frac{3}{2}x. \quad (85)$$

Полиномы Лежандра обладают следующими свойствами:

1. Полином $L_{n+1}(x)$ номера $n+1$ является полиномом n -ой степени, обладающим той же четностью, что и n .

$$P_n(-x) = (-1)^n P_n(x). \quad (86)$$

2. Полиномы Лежандра $P_n(x)$ в точках $x = \pm 1$ принимают следующие значения:

$$P_n(1) = 1, \quad P_n(-1) = (-1)^n. \quad (87)$$

Эта формула, как и формулы (32), (33) для методов прямоугольников и трапеций, выделяет в остаточном члене γ_n главное слагаемое $C/n!$, которое стремится к нулю как $n \rightarrow \infty$. Коэффициент C (57) не зависит от n . Дополнительное слагаемое $\sigma_n/n!$ является бесконечно малой более высокого порядка. Если им пренебречь, то получится асимптотическое представление остаточного члена

$$\gamma_n \approx Cn!. \quad (59)$$

Его относительная точность возрастает с увеличением n .

Другое представление остаточного члена γ_n можно вывести с помощью формулы (39). Она позволяет записать формулу (53) в виде

$$\gamma_n = \frac{(b-a)^3}{180n^4} f^{(4)}(\eta), \quad (60)$$

где $\eta \in [x_{j-2}, x_j]$. Существование такой точки гарантировано, но ее точное положение на отрезке неизвестно.

Суммируя равенства (52) по j , получим квадратурную формулу (18) со следующим выражением для остаточного члена:

$$|\gamma_n| \leq \frac{h^3}{90} \sum_{j=1}^{n/2} f^{(4)}(\eta_j). \quad (61)$$

Данная оценка позволяет определить, с каким n нужно проводить вычисления, чтобы погрешность не превышала заданной точности ε . Кроме того, если четвертая производная функции $f(x)$ при любом значении аргумента x на отрезке $[a, b]$ и все в оценки же (43), (44), (61) входят константы M_2 и M_4 , мажущие вторую и четвертую производные в граничных точках отрезка $[a, b]$. Тогда

$$|\gamma_n| \leq \frac{(b-a)^3}{180n^4} M_4. \quad (62)$$

Чтобы определить, с каким n нужно проводить вычисления, чтобы погрешность не превышала заданной точности ε , нужно решить неравенство

Метод Симпсона является методом более высокого порядка точности – четвертого. В этом его преимущество перед методами прямоугольников и трапеций. Правда, приведенные выше оценки остаточного члена, требуют большей гладкости подинтегральной функции – она должна быть четыре раза непрерывно дифференцируема.

Задача 2. Вычислить интеграл (46) по формуле Симпсона при $n=2$.

В данном случае

$$S_2 = \frac{\pi}{12} \left(\sin 0 + 4 \sin \frac{\pi}{4} + \sin \frac{\pi}{2} \right) = 1.002280, \quad (63)$$

$$\gamma_2 = -0.002280. \quad (64)$$

Четвертая производная функции $\sin x$ на отрезке $\left[0, \frac{\pi}{2}\right]$ положительна и не превосходит единицы, так что знак погрешности согласуется с формулой (60), а ее величина – с оценкой (61):

$$|\gamma_2| \leq \frac{1}{180 \cdot 16} \left(\frac{\pi}{2} \right)^3 < 0.0034. \quad (65)$$

$$I = \frac{16}{15} S_2 - \frac{1}{15} S_{n/2} + \tilde{\gamma}_n, \quad (66)$$

где $\tilde{\alpha}_n$, $\tilde{\beta}_n$, $\tilde{\gamma}_n$ – остаточные члены этих модифицированных формул

$$\tilde{\alpha}_n = \frac{4}{3n!} (\mu_n - \mu_{n/2}) = o(n^{-1}), \quad (67)$$

$$\tilde{\beta}_n = \frac{4}{3n!} (\nu_n - \nu_{n/2}) = o(n^{-1}), \quad (68)$$

$$\tilde{\gamma}_n = \frac{16}{15n!} (\sigma_n - \sigma_{n/2}) = o(n^{-1}). \quad (69)$$

Формулы (71), (72), (73), написанные по результатам двух расчетов с числом точек $n/2$ и n , являются асимптотически более точными, чем исходные. В исходных формулах погрешности убывают, соответственно, как n^{-2} , n^{-2} , n^{-4} , в модифицированных формулах погрешности согласно (74), (75), (76) являются бесконечно малыми более высокого порядка. Однако для исходных формул известны оценки погрешностей (43), (44), (61). Для модифицированных формул в нашем распоряжении оценок нет. Если мы хотим ими пользоваться, то нужно провести соответствующее исследование. Исключение составляет формула (72). Согласно формуле (21) ее можно переписать в виде

$$I = S_n + \tilde{\beta}_n, \quad (70)$$

т. е. модифицированная формула трапеций оказывается просто формулой Симпсона с уже известным остаточным членом $\gamma_n = \tilde{\beta}_n$.

Задача 3.

Вычислить по формуле Симпсона интеграл (46) с $n=4$. Используя результаты задачи 2, найти приближенную апостериорную погрешность (70).

В данном случае

$$S_4 = \frac{24}{24} \left(\sin 0 + 4 \sin \frac{\pi}{8} + 2 \sin \frac{\pi}{4} + 4 \sin \frac{3\pi}{8} + \sin \frac{\pi}{2} \right) = 1.000135, \quad (71)$$

$$\gamma_4 = -0.000135. \quad (72)$$

Апостериорная оценка погрешности по результатам двух расчетов дает

$$\gamma_4 \approx \frac{1}{15} (S_4 - S_2) = -0.000143. \quad (73)$$

Несмотря на маленкое число точек, она хорошо согласуется с фактической погрешностью (78), соединенной «в лоб» по известному значению интеграла (46).

Задача 4.

Используя результаты решения задач 2 и 3, посчитать интеграл (46) по модифицированной формуле Симпсона (73).

В данном случае

$$I = \frac{16}{15} S_4 - \frac{1}{15} S_{n/2} + \tilde{\gamma}_4, \quad (74)$$

где $\tilde{\alpha}_4$, $\tilde{\beta}_4$ – остаточные члены в формуле (73).

На практике этого параграфа каждый раз, когда мы будем говорить о производных полиномов какой-нибудь степени, всегда будем включать в них полиномы более низких степеней, не оговаривая это особо.

Полагая последовательно $f(x) = 1, x, x^2, \dots, x^{2n-1}$ и принимая во внимание, что для этих функций, согласно требованию Гаусса, остаточный член должен равняться нулю, получим:

$$\int x^n dx = \frac{1}{(m+1)!} (1 - (-1)^m) = \sum_{i=1}^m c_i x^i, \quad 0 \leq m \leq 2n-1. \quad (80)$$

Соотношения (81) представляют собой систему $2n$ нелинейных уравнений с $2n$ неизвестными, в качестве которых выступают узлы x_i и веса c_i ($1 \leq i \leq n$).

Уравнение (81), соответствующее индексу $m=0$, дает

$$\sum_{i=1}^m c_i = 2. \quad (81)$$

Полагая последовательно $f(x) = 1, x, x^2, \dots, x^{2n-1}$ и принимая во внимание, что для этих функций, согласно требованию Гаусса, остаточный член должен равняться нулю, получим:

$$J = (-1)^{m+1} \frac{1}{2^n n!} \int_{-1}^1 \frac{d^{m+1} Q_m(x)}{dx^{m+1}} \frac{d^{m+1}}{dx^{m+1}} (x^2 - 1)^m dx = 0. \quad (82)$$

Здесь под знаком интеграла в качестве множителя стоит $(m+1)$ -я производная полинома m -ой степени $Q_m(x)$, тождественно равная нулю. Ортогональность доказана.

Сделаем важное замечание. Соотношение ортогональности (87) справедливо, в частности, в случае, когда в качестве полинома $Q_m(x)$ взят полином Лежандра $P_m(x)$:

$$\int_{-1}^1 P_m(x) P_n(x) dx = 0, \quad \text{при } m < n. \quad (83)$$

Фактически в этом условии ортогональности не важно, какой именно из двух индексов m или n больше, а какой меньше. Важно лишь, что они не равны. Таким образом, из свойства 4 вытекает следствие.

Следствие 1.

Полиномы Лежандра образуют систему полиномов, ортогональных на отрезке $[-1, 1]$:

$$\int_{-1}^1 P_m(x) P_n(x) dx = 0, \quad \text{при } m \neq n. \quad (84)$$

Из линейной алгебры известно, что система полиномов, ортогональных на некотором множестве, определяется однозначно с точностью до множителей. Поэтому следствию 1 можно поставить обратное утверждение.

Следствие 2.

Любая система полиномов, ортогональных на отрезке $[-1, 1]$, совпадает с множеством полиномов Лежандра.

3.3. Узлы и весовые коэффициенты квадратурных формул Гаусса.

Изучив свойства полиномов Лежандра, передадим к решению основной задачи – определению узлов и весовых коэффициентов квадратурных формул Гаусса. Составим полином n -ой степени

$$\phi_n(x) = (x - x_1)(x - x_2) \cdots (x - x_n), \quad (90)$$

где x_i – искомые узлы. Возьмем произвольный полином $Q_m(x)$ степени $m < n$, помножим его на полином $\phi_n(x)$ и проинтегрируем произведение по отрезку $[-1, 1]$ с помощью квадратурной формулы (80). Поскольку это произведение представляет

собой полином степени $m+n \leq 2n-1$, формула Гаусса должна быть для него точной. В результате согласно (90) получим:

$$\int_{-1}^1 Q_m(x) \phi_n(x) dx = \sum_{i=1}^n c_i Q_{n-i}(x) \phi_n(x) = 0. \quad (91)$$

Мы видим, что полином $\phi_n(x)$ ортогонален к любому полиному степени $m < n$ в том числе и к полиномам Лежандра $m < n$. Это означает, что он с точностью до мономов совпадает с n -ым полиномом Лежандра $\phi_n(x) = A_n P_n(x)$. Отсюда следует вывод: узлы квадратурной формулы Гаусса являются корнями полинома Лежандра $P_n(x)$. Напомним, что корни полиномов Лежандра расположены на интервале $(-1, 1)$ симметрично относительно ее средней точки $x = 0$.

Для того, чтобы подсчитать весовые коэффициенты c_i , введем специальные полиномы

$$Q_{n-1,m}(x) = \frac{(x-x_1)\cdots(x-x_{m-1})(x-x_{m+1})\cdots(x-x_n)}{(x_n-x_1)\cdots(x_n-x_{m-1})(x_n-x_{m+1})\cdots(x_n-x_n)}. \quad (92)$$

Каждый из них является полиномом степени $(n-1)$. В числителе у него стоит полином $\phi_n(x)$ с опущенным множителем $(x-x_n)$, знаменатель - значение числителя в точке $x = x_n$. В результате такой структуры полином $Q_{n-1,m}(x)$ в точках x_i удовлетворяет соотношениям:

$$Q_{n-1,m}(x_i) = \begin{cases} 0, & i \neq m \\ 1, & i = m \end{cases}. \quad (93)$$

Для полинома $Q_{n-1,m}(x)$ квадратурная формула Гаусса должна быть точной. С учетом (93) это дает

$$\int_{-1}^1 Q_{n-1,m}(x) dx = \sum_{i=1}^n c_i Q_{n-1,m}(x_i) = c_m. \quad (94)$$

В результате получаем следующее интегральное выражение для весовых коэффициентов квадратурной формулы Гаусса:

$$c_m = \int_{-1}^1 Q_{n-1,m}(x) dx = \int_{-1}^1 \frac{(x-x_1)\cdots(x-x_{m-1})(x-x_{m+1})\cdots(x-x_n)}{(x_n-x_1)\cdots(x_n-x_{m-1})(x_n-x_{m+1})\cdots(x_n-x_n)} dx. \quad (95)$$

3.4. Исследование квадратурной формулы.

Нам остались решить последние вопросы - доказать, что квадратурная формула, у которой в качестве узлов x_i берется корни полинома Лежандра, а весовые коэффициенты c_i вычисляются по формуле (95), действительно решает задачу Гаусса, являясь точной для любого полинома степени $(2n-1)$.

Проведем доказательство в два этапа. Сначала докажем, что такая формула является точной для любого полинома $Q_{n-1,m}(x)$ степени $(n-1)$. Такой полином можно представить в виде суммы специальных полиномов (92).

Формулы Симпсона и Гаусса дают в данном случае следующие результаты:

$$S_2 = \frac{1}{3}(e^{-1} + 4 + e) = \frac{4}{3} \cdot \frac{2}{3} chl = 2.362054,$$

$$\gamma_2 = -0.011651,$$

$$G_2 = (e^{-1/\sqrt{5}} + e^{1/\sqrt{5}}) = 2ch\frac{1}{\sqrt{5}} = 2.342696,$$

$$\delta_2 = 0.007706.$$

Мы видим, что даже с двумя узлами формула Гаусса дает хороший ответ. Его точность выше точности ответа, полученного по формуле Симпсона.

В заключение сделаем следующее замечание. Несмотря на высокую точность квадратурных формул Гаусса, при компьютерных расчетах ими пользуются сравнительно редко. Дело в том, что для применения метода Гаусса нужно либо ввести в компьютере до начала расчетов корни полинома Лежандра и весовые коэффициенты, либо составить специальную подпрограмму для их вычисления. В результате потеря человеческого и машинного времени на подготовку программы к основному расчету, связанным с вычислением интеграла, могут не окупиться точностью метода Гаусса. Вычисление интеграла по более простой схеме метода Симпсона имеет от этой точки зрения преимущество.

§4. Построение первообразной с помощью численного интегрирования

Формулы Ньютона-Лейбница (1) позволяют выразить значения определенного интеграла от функции $f(x)$ через ее первообразную $F(x)$. В математическом анализе устанавливается и прямо противоположная возможность: первообразная функции $f(x)$, непрерывной на отрезке $[a, b]$, может быть записана в виде определенного интеграла с переменным верхним пределом:

$$F(x) = \int_a^x f(t) dt. \quad (106)$$

Здесь x_0, x - две точки отрезка $[a, b]$, причем нижний предел интегрирования x_0 предполагается фиксированным, верхний x - переменным. В случае непрерывной функции $f(x)$ функция $F(x)$, определенная с помощью интеграла (106), является дифференцируемой и ее производная равна $f(x)$:

$$F'(x) = \frac{d}{dx} \left(\int_a^x f(t) dt \right) = f(x). \quad (107)$$

Формула (106) в сочетании с какой-нибудь формулой численного интегрирования, например, Симпсона, представляет собой универсальный алгоритм построения первообразной. Приведем два примера, иллюстрирующие этот алгоритм.

Функция $f(x) = \sin x$ непрерывна и, следовательно, имеет первообразные. Они не могут быть выражены через элементарные функции, но представление в виде интеграла с переменным верхним пределом для них справедливо. Одну из

$$Q_{n-1}(x) = \sum_{i=1}^n Q_{n-1}(x_i) Q_{n-1,n}(x). \quad (96)$$

Справедливость данного разложения вытекает из следующих соображений. Здесь левая и правая части равенства совпадают в n точках $x_i, 1 \leq i \leq n$. Но, если два полинома $(n-1)$ -й степени совпадают в n точках, то они тождественно равны.

Интегрируя равенство (96) по отрезку $[-1, 1]$, получим

$$\int_{-1}^1 Q_{n-1}(x) dx = \sum_{i=1}^n Q_{n-1}(x_i) \int_{-1}^1 Q_{n-1,n}(x) dx = \sum_{i=1}^n c_i Q_{n-1,n}(x_i). \quad (97)$$

Итак, для полиномов $(n-1)$ -й степени утверждение доказано.

Теперь рассмотрим произвольный полином $Q_{n,m}(x)$ степени $(2n-1)$. Разделим его с остатком на полином Лежандра $P_n(x)$ и представим в виде:

$$Q_{n-1,m}(x) = P_n(x) q_{n-1}(x) + r_{n-1}(x). \quad (98)$$

где $q_{n-1}(x)$ и $r_{n-1}(x)$ полиномы степени $(n-1)$. Пронизитрировав равенство (98) по отрезку $[-1, 1]$, будем иметь:

$$\begin{aligned} \int_{-1}^1 Q_{n-1,m}(x) dx &= \int_{-1}^1 \{P_n(x) q_{n-1}(x) + r_{n-1}(x)\} dx = \int_{-1}^1 r_{n-1}(x) dx = \\ &= \sum_{i=1}^n c_i r_{n-1}(x_i) = \sum_{i=1}^n c_i \{P_n(x_i) q_{n-1}(x_i) + r_{n-1}(x_i)\} = \sum_{i=1}^n c_i Q_{n-1,n}(x_i) \end{aligned} \quad (99)$$

Появился выполненный преобразование. Интеграл $\int_{-1}^1 P_n(x) q_{n-1}(x) dx$ опущен, поскольку

полином Лежандра $P_n(x)$ ортогонален к любому полиному $(n-1)$ -й степени. Оставшийся интеграл от полинома $r_{n-1}(x)$ вычислен с помощью квадратурной формулы (97). Выше уже доказано, что для полиномов степени $(n-1)$ она является точной.

Последний переход заключается в том, что в сумме $\sum_{i=1}^n c_i r_{n-1}(x_i)$ добавлены слагаемые $P_n(x_i) q_{n-1}(x_i)$. Они не меняют значения суммы, поскольку все равны нулю: ведь узлами квадратурной формулы являются корни полинома Лежандра $P_n(x)$.

Итак, построенная квадратурная формула действительно является точной для любого полинома степени $(2n-1)$, т. е. задача Гаусса решена. На основе погрешности квадратурных формул Гаусса мы останавливаться не будем, однако задачи, к разбору которых переходим, показывают, что эти формулы обеспечивают для гладких функций очень высокую точность.

Задача 5.

Построить квадратурную формулу Гаусса с двумя и тремя узлами.

Выведем сначала квадратурную формулу с двумя узлами. Узлы определяются как корни второго полинома Лежандра, выражение для которого мы выписали выше (85). В данном случае имеем:

$$x_1 = -1/\sqrt{3}, x_2 = 1/\sqrt{3}.$$

Узлы расположены симметрично относительно точки $x = 0$.

Бесовские коэффициенты расчитываются по формуле (95):

$$c_1 = \frac{1}{2} \frac{x - x_2}{x_2 - x_1} dx = \frac{2 - 1/\sqrt{3}}{-2/\sqrt{3}} dx = 1,$$

$$c_2 = \frac{1}{2} \frac{x - x_1}{x_2 - x_1} dx = \frac{2 + 1/\sqrt{3}}{2/\sqrt{3}} dx = 1.$$

Они равны между собой, а их сумма, в соответствии с общим соотношением (82), равна двум. В результате имеем квадратурная формула принимает вид:

$$\int_{-1}^1 f(x) dx = f(-1/\sqrt{3}) + f(1/\sqrt{3}) + \delta_1. \quad (102)$$

Она является точной для любого полинома третьей степени.

Перейдем теперь к выводу квадратурной формулы Гаусса с тремя узлами.

Согласно формуле (85) для третьего полинома Лежандра ее узлами являются числа:

$$x_1 = -\sqrt{3}/5, x_2 = 0, x_3 = \sqrt{3}/5.$$

Остается подсчитать весовые коэффициенты:

$$c_1 = \frac{1}{3} \frac{x - x_2}{x_3 - x_1} dx = \frac{x + \sqrt{3}/5}{(\sqrt{3}/5)(-2\sqrt{3}/5)} dx = \frac{5}{9},$$

$$c_2 = \frac{1}{3} \frac{x - x_1}{x_3 - x_2} dx = \frac{x + \sqrt{3}/5}{(\sqrt{3}/5)(2\sqrt{3}/5)} dx = \frac{5}{9},$$

В результате квадратурная формула Гаусса с тремя узлами запишется в виде:

$$\int_{-1}^1 f(x) dx = \frac{5}{9} f(-\sqrt{3}/5) + \frac{8}{9} f(0) + \frac{5}{9} f(\sqrt{3}/5) + \delta_2. \quad (105)$$

Она является точной для любого полинома пятой степени.

Задача 6.

Вычислить по формулам Симпсона и Гаусса при $n = 2$ интеграл:

$$\int_{-1}^1 e^x dx = 2shl = 2.350402.$$

Сравнить результаты численного интегрирования с точным значением интеграла и между собой.

первообразных мы получим, выбирая нижний предел интегрирования $x_0 = 0$. Ее называют интегральным синусом и обозначают

$$Si(x) = \int_0^x \frac{\sin t}{t} dt.$$

Интегральный синус определен на всей числовой прямой, является нечетной функцией x , имеет конечные предельные значения на бесконечности

$$\lim_{x \rightarrow \pm\infty} Si(x) = \pm \frac{\pi}{2}.$$

Согласно (107)

$$Si'(x) = \sin x/x.$$

По знаку производной легко определить области возрастания и убывания функции, разделенные точками экстремума $x_k = k\pi$ ($k = \pm 1, \pm 2, \dots$). Методы численного интегрирования позволяют вычислить значения $Si(x)$ при любом x . График интегрального синуса при $x \geq 0$ приведен на рис. 3.

В качестве второго примера рассмотрим функцию ошибок $erf(x)$, играющую важную роль в теории вероятности. Ее обозначение образовано с помощью первых букв английского названия функции ошибок - error function. Подобно интегральному синусу, функция ошибок вводится в виде интеграла с переменным пределом от функции e^{-x^2} , которая не имеет первообразных в классе элементарных функций:

$$erf(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt.$$

Функция ошибок определена на всей числовой прямой, является нечетной функцией x , имеет конечные предельные значения на бесконечности:

$$\lim_{x \rightarrow \pm\infty} erf(x) = \pm 1.$$

Согласно (107)

$$erf'(x) = \frac{2}{\sqrt{\pi}} e^{-x^2}.$$

Производная всюду положительна, следовательно, функция ошибок монотонно возрастает. Ее график приведен на рис. 4.

Существует ряд других специальных функций, которые вводятся как интегралы с переменным верхним пределом. Не будем останавливаться на их описание, отметим лишь, что разбранные примеры показывают, насколько условно деление функций на элементарные и неэлементарные. По существу, чтобы работать с какой-нибудь функцией, нужно знать ее свойства и иметь алгоритм вычисления при любом значении аргумента. С этой точки зрения применение интегрального синуса или функции ошибок ничем не отличается от применения привычных нам элементарных функций.

Это определение законно, поскольку узловые охватывает тремя аксиомам нормы:

$$\|x\| \geq 0,$$

причем равнота нулю имеет место только для нулевого элемента.

2. Модуль числового множества можно вычислить за знак нормы

$$\|xy\| = |z| \|y\|.$$

3. Неравенство треугольника

$$\|x+y\| \leq \|x\| + \|y\|.$$

Справедливость последнего утверждения вытекает из свойства максимума:

$$\max_{0 \leq x \leq a} |y + z| \leq \max_{0 \leq x \leq a} |y| + \max_{0 \leq x \leq a} |z|.$$

1.2. РАЗНОСТНЫЕ АППРОКСИМАЦИИ ПЕРВОЙ ПРОИЗВОДНОЙ.

Для сеточных функций нельзя ввести обычное понятие производной, включающее операцию предельного перехода при $\Delta x \rightarrow 0$. Вместо производной здесь вводятся разностные отношения:

$$\|y\|^{(1)} = \max_{0 \leq x \leq a} |y(x)|.$$

Глава 4. §1.

Перейдем к априорной оценке погрешности. Вторая и третья производные рассматриваемой функции $y(x)$ имеют вид

$$y''(x) = \frac{2}{(1-x)^3}, \quad y'''(x) = \frac{6}{(1-x)^4}.$$

Для них на отрезке $[-0.1, 0.1]$ справедливы оценки

$$|y''(x)| \leq \frac{2}{(0.9)^3} < 2.8, \quad |y'''(x)| = \frac{6}{(0.9)^4} \leq 9.3.$$

Так что неравенства (15) и (19) записываются следующим образом

$$\|y'\| \leq 0.14, \quad \|y''\| \leq 0.14, \quad \|y'''(x)\| \leq 0.14, \quad \|y^{(4)}(x)\| \leq 0.031.$$

Они выполняются.

1.3. РАЗНОСТНАЯ АППРОКСИМАЦИЯ ВТОРОЙ ПРОИЗВОДНОЙ.

Для разностной аппроксимации второй производной составим разностное отношение первых разностных производных

$$\frac{y_{i+1} - y_i}{h} - \frac{y_i - y_{i-1}}{h} = \frac{h}{2} (y_{i+1} - y_i) - \frac{h}{2} (y_i - y_{i-1}).$$

Чтобы установить связь выражения (21) со второй производной, предположим, что на отрезке $[a, b]$ определены дифференцируемые функции $y(x)$, значения которых в точках сетки x_i дают значения сеточной функции y_i . Вычислим ее вторую производную в тех же точках x_i , и составим разности:

$$y''_i = L_0[y_i] - y''(x_i), \quad 1 \leq i \leq n-1. \quad (22)$$

Она представляет собой погрешность аппроксимации второй производной с помощью разностного отношения второго порядка (21).

Оценим величину погрешности при предположении, что функция $y(x)$ четырех раз непрерывно дифференцируема на отрезке $[a, b]$. Это предположение позволяет написать разложение Тейлора

$$y_{i+1} = y(x_i + h) = y_i + y'(x_i)h + \frac{1}{2} y''(x_i)h^2 + \frac{1}{6} y'''(x_i)h^3 + \frac{1}{24} y^{(4)}(x_i)h^4,$$

$$y_{i-1} = y(x_i - h) = y_i - y'(x_i)h + \frac{1}{2} y''(x_i)h^2 - \frac{1}{6} y'''(x_i)h^3 + \frac{1}{24} y^{(4)}(x_i - h)h^4.$$

Подставляя их в формулу (21), получим

$$y''_i = \frac{1}{24} (y^{(4)}(x_i + h) + y^{(4)}(x_i - h)) h^2. \quad (24)$$

Мы не можем вычислить погрешность по этой формуле, поскольку значения аргументов $y^{(4)}(x)$ нам неизвестны, но можем ее оценить. Функция $y^{(4)}(x)$ непрерывна и, следовательно, ограничена на отрезке

$$|y^{(4)}(x)| \leq M_4, \quad a \leq x \leq b.$$

В результате из формулы (24) получаем

$$\|\psi\| \leq \frac{1}{12} M_1 h^2. \quad (26)$$

Таким образом, разностное отношение (21) аппроксимирует вторую производную со вторым порядком точности относительно h для функций, имеющих четыре непрерывные производные на отрезке $[a, b]$. Совершенно аналогично можно строить разностные аналоги производных более высокого порядка.

Задача 2.

Для функции $y = 1/(1-x)$ вычислить на сетке (20) вторую разностную производную в точке $x_1 = 0$. Найти погрешность аппроксимации второй производной $y''(0) = 2$ и сравнить результат с априорной оценкой (26).

В данном случае

$$L_h[\psi] = \frac{1/1.1 - 2 + 1/0.9}{0.01} = 2.020202, \\ u_i = 0.020202.$$

Четвертая производная рассматриваемой функции $y(x)$ и мажоранта для нее на отрезке $[-0.1, 0.1]$ имеют вид

$$y^{(4)}(x) = \frac{24}{(1-x)^5}, \quad |\psi^{(4)}(x)| \leq \frac{24}{(0.9)^5} < 41.$$

Так что неравенство (26) записывается следующим образом

$$|\psi| \leq 0.034.$$

Оно выполняется.

При численном интегрировании дифференциальных уравнений производные в них приближенно заменяются соответствующими разностными отношениями. В результате задача сводится к системе разностных уравнений, которые решаются на компьютере. В качестве ответа получается сеточная функция $\{\psi_i\}$ ($0 \leq i \leq n$). После этого встает вопрос, в какой степени и с какой точностью ее можно рассматривать в качестве приближенного решения исходной задачи. Нужно иметь в виду, что прямое сравнение решения дифференциального уравнения $u(x)$ и рассчитанной сеточной функции невозможно: они принадлежат разным пространствам и их, прежде всего, нужно свести в одно пространство. Это можно сделать двумя способами.

В первых, по сеточной функции с помощью методов интерполяции можно построить функцию непрерывного аргумента $u(x)$ и оценить разность $z(x) = u(x) - u(x)$, например, в норме C

$$\|z\| = \max_i |z_i| = \|u(x) - u(x)\|.$$

Во вторых, наоборот, решение дифференциального уравнения можно сопоставить сеточную функцию $\{u_i = u(x_i)\}$ и сравнить между собой две сеточные функции $\{\psi_i\}$ и

Введем для оценки сеточной функции ψ ее норму

$$\|\psi\|_{\text{норм}} = \max_i |\psi_i|, \text{ при этом } \|\psi\| \leq \|\psi\|_{\text{норм}}. \quad (39)$$

Предположим далее, что функция $\frac{d}{du}(x, u)$ в интересующей нас области изменения ее аргументов ограничена

$$\left| \frac{d}{du}(x, u) \right| \leq C.$$

Это позволяет написать оценку

$$\left| 1 + h \frac{d}{du}(x, u + \theta h) \right| \leq 1 + Ch < e^{Ch} = q, \quad q > 1. \quad (41)$$

С учетом (39) и (41) из формулы (37) следуют рекуррентные неравенства

$$|z_i| \leq |z_{i-1}| + \|\psi\| h, \quad (42)$$

которые порождают цепочку оценок

$$\begin{aligned} z_0 &= 0, \\ |z_i| &\leq \|\psi\| h, \\ |z_i| &\leq (1+q) \|\psi\| h, \\ |z_i| &\leq (1+q+q^2) \|\psi\| h, \\ &\vdots \\ |z_i| &\leq (1+q+q^2+\dots+q^{i-1}) \|\psi\| h. \end{aligned} \quad (43)$$

Согласно (41) $q > 1$, так что

$$1 + q + q^2 + \dots + q^{i-1} < nq^i = ne^{Ci}. \quad (44)$$

Это позволяет заменить индивидуальные оценки (43) универсальной оценкой

$$|z_i| \leq nh^{Cm} \|\psi\|, \quad 0 \leq i \leq n. \quad (44)$$

Неравенства (44) справедливы при любом i , в частности, при том, при котором $|z_i|$ достигает своего наибольшего значения и определяет тем самым норму сеточной функции $\|\psi\|$. В результате оценка погрешности решения принимает вид

$$\|z\| \leq h e^{Cm} \|\psi\|, \quad (45)$$

где l – длина отрезка, на котором рассматривается решение исходной задачи (27), (28).

Мы получили важный результат: оценку погрешности решения через оценку погрешности аппроксимации уравнения с коэффициентом, который не зависит от шага h . Чем лучше разностное уравнение аппроксимирует дифференциальное, тем меньше погрешность решения.

Чтобы завершить исследование метода Эйлера, оценим норму погрешности аппроксимации уравнения $\|\psi\|$. Предположим, что функция $f(x, u)$ имеет в рассматриваемой области изменения аргументов непрерывные и ограниченные первые

Предположим, что решение дифференциального уравнения $u(x)$ имеет производные достаточно высокого порядка и напишем для него разложение по формуле Тейлора

$$u_{i+1} = u_i + u'_i(x)h + \frac{1}{2}u''_i(x)h^2 + \frac{1}{6}u'''_i(x)h^3 + \dots \quad (54)$$

Если его обозвать на члене порядка h и положить в соответствия с дифференциальным уравнением (27) $y_i = f(x_i, u_i)$, то мы приедем к схеме Эйлера.

Сделаем следующий шаг. Оборнем разложение (54) на члене порядка h^2 и воспользуемся для вычисления производной $u'(x)$ формулой (46). В результате получим новое рекуррентное соотношение, более сложное чем (32),

$$y_{i+1} = y_i + f(x_i, y_i)h + \frac{1}{2} \left[\frac{d}{dx}(x_i, y_i) + \frac{d}{dy}(x_i, y_i)f(x_i, y_i) \right] h^2, \quad (55)$$

которое можно также записать в виде разностного уравнения

$$\frac{y_{i+1} - y_i}{h} = y_i + f(x_i, y_i) + \frac{1}{2} \left[\frac{d}{dx}(x_i, y_i) + \frac{d}{dy}(x_i, y_i)f(x_i, y_i) \right] h. \quad (56)$$

Далее, как и в предыдущем разделе, мы обозначили исключенную в разностном уравнении (56) буквой y , а не u , чтобы подчеркнуть, что (27) и (56) – это два разных уравнения.

Уравнение (55), дополненное начальными условиями (31), дает явную разностную схему численного решения рассматриваемой задачи Коши. По рекуррентной формуле можно последовательно рассчитать все значения сеточной функции y_i , $0 \leq i \leq n$ и получить таким образом приближенное решение задачи (27), (28). Исследование показывает, что такая усложненная схема имеет второй порядок точности относительно h как для аппроксимации уравнения, так и для погрешности решения. Существенно то, что основная идея данного подхода допускает дальнейшее развитие.

Если оборвать разложение (54) на члене порядка h^3 , h^4 и т. д., то получатся разностные схемы третьего, четвертого и более высоких порядков точности.

Однако у данного подхода есть существенный недостаток. При расчетах по схеме Эйлера требуется вычислять только значения функции $f(x, y)$. В схеме же (55) на каждом шаге приходится вычислять не только функцию f , но и ее первые производные $\frac{d}{dx}(x, y)$ и $\frac{d}{dy}(x, y)$. Если мы, оставив в разложении (54) члены до h^4 включительно, построим схему четвертого порядка точности, то на каждом шаге придется вычислять десять величин: функцию $f(x, y)$, две ее первых производных, три вторых производных и четыре третьих производных. Это существенно усложняет разработку программы и нарушит важный принцип вычислительной математики – использовать в расчетах только те величины, которые заданы условиями задачи. Формулировка задачи Коши предполагает, что известен алгоритм вычисления функции $f(x, u)$ по значениям ее аргументов. Если этот алгоритм сводится к расчету по простой формуле, то вычисление производных не составляет труда. Однако

$\{y_i\}$, составив их разность $z_i = y_i - u_i$. При этом погрешность приближенного решения задачи будет характеризовать норма разности

$$\|z\| = \max_i |y_i - u_i|.$$

Наиболее последовательным является первый путь, но обычно выбирают более простой – второй.

В следующих параграфах мы рассмотрим численное решение с помощью метода конечных разностей задачи Коши и краевой задачи для линейного дифференциального уравнения второго порядка.

9.2. Численное решение задачи Коши.

Рассмотрим задачу Коши для дифференциального уравнения первого порядка

$$u' = f(x, u),$$

$$u(x_0) = u_0.$$

Если функция $f(x, u)$ непрерывна и удовлетворяет условию Липшица по аргументу x в некоторой окрестности начальной точки (x_0, u_0) , то можно указать такой отрезок $[a, b]$, $a < x_0 < b$, на котором решение задачи (27), (28) $u(x)$ существует и является единственным. В этом параграфе мы обсудим численные методы ее решения.

2.1. Метод Эйлера.

Пусть нам нужно построить решение задачи (27), (28) на отрезке $[x_0, x_0 + l]$. Возьмем некоторое целое число n , введем шаг $h = l/n$ и образуем на отрезке сетку

$$x_i = x_0 + ih, \quad 0 \leq i \leq n.$$

Сопоставим задаче (27), (28) на отрезке разностную задачу

$$\frac{y_{i+1} - y_i}{h} = f(x_i, y_i), \quad 0 \leq i \leq n-1; \quad (29)$$

$$y_0 = u_0.$$

Здесь мы заменили производную $u'(x)$ в уравнении (27) правой разностной производной и сохранили неизменным начальное условие (28).

Уравнение (29) является разностным уравнением первого порядка, которое принято называть схемой Эйлера. Его можно переписать в виде рекуррентного соотношения

$$y_{i+1} = y_i + h f(x_i, y_i), \quad 0 \leq i \leq n-1. \quad (30)$$

Это позволяет последовательно рассчитать все значения сеточной функции $\{y_i\}$, решив тем самым задачу (29), (30). Такую разностную схему называют явной.

Перейдем теперь к обсуждению главного вопроса: с какой точностью рассчитанная сеточная функция $\{y_i\}$ есть решение исходной задачи Коши $u(x)$. Для ответа на него рассмотрим решение задачи (27), (28) в сетках (29), образуя из функции непрерывного аргумента сеточную функцию $\{u_i = u(x_i)\}$, и сравнив ее с

рассчитанной сеточной функцией $\{y_i\}$. Для этого образуем две сеточные функции \mathbf{z} ,

$$\mathbf{z} = y_i - u_i, \quad 0 \leq i \leq n;$$

$$\psi_i = \frac{u_{i+1} - u_i}{h} - f(x_i, u_i), \quad 0 \leq i \leq n-1. \quad (34)$$

Смысл первой функции (33) очевиден. Она характеризует разницу между рассчитанными числами y_i и u_i , решением $u(x)$ задачи (27), (28) в точках сетки x_i . В соответствии с этим сеточную функцию \mathbf{z} называют погрешностью решения.

Функция ψ (34) получается в результате подстановки решения дифференциального уравнения (27) в разностное уравнение (30). Если бы эти уравнения совпадали, то мы получили бы нуль. Но они различаются и нуль мы не получим. Сеточную функцию ψ , характеризующую степень близости дифференциального и разностного уравнений, называют погрешностью аппроксимации уравнения на решении.

Установим связь между сеточными функциями \mathbf{z} и ψ . С этой целью выразим из формулы (33) y_i :

$$y_i = u_i + \mathbf{z}_i, \quad (35)$$

и подставим в разностное уравнение (30). В результате получим

$$\frac{y_{i+1} - y_i}{h} + \frac{u_{i+1} - u_i}{h} - f(x_i, u_i + \mathbf{z}_i) =$$

или

$$\frac{y_{i+1} - y_i}{h} = \left\{ f(x_i, u_i + \mathbf{z}_i) - f(x_i, u_i) \right\} - \frac{u_{i+1} - u_i}{h} - f(x_i, u_i). \quad (36)$$

Здесь в обе фигуры скобки мы добавили величину $f(x_i, u_i)$. Добавленные члены входят в соотношение (36) с противоположными знаками и благодаря этому не нарушают равенство. После таких преобразований во вторых фигурах скобках получается величина ψ_i .

В первых фигурах скобках стоит разность значений функции f при одинаковом первом аргументе x_i и разных значениях второго аргумента. Этую разность с помощью формулы Лагранжа можно представить в виде

$$f(x_i, u_i + \mathbf{z}_i) - f(x_i, u_i) = \frac{d}{du} f(x_i, u_i + \theta_i \mathbf{z}_i),$$

и записать формулу (36) в виде рекуррентного соотношения

$$z_{i+1} = \left\{ 1 + \frac{d}{du} f(x_i, u_i + \theta_i \mathbf{z}_i) \right\} z_i - \psi_i, \quad 0 \leq i \leq n-1. \quad (37)$$

Согласно (28) и (31) ему следует дополнить нулем начальным условием

$$z_0 = 0. \quad (38)$$

В отличии от формулы (30), (31) формулы (37), (38) не могут быть использованы для вычисления величин ψ_i . В них входят неизвестные величины ψ_i , u_i , θ_i . Однако из этой системы рекуррентных равенств можно получить рекуррентные неравенства.

частные производные $\frac{\partial f}{\partial x}$ и $\frac{\partial f}{\partial u}$. Это обеспечивает существование у решения задачи (27), (28) непрерывной и ограниченной второй производной

$$u''(x) = \frac{\partial f}{\partial x}(x, u) + \frac{\partial f}{\partial u}(x, u) f(x, u). \quad (46)$$

Запишем для функции $u(x)$ формулу Тейлора с остаточным членом в форме Лагранжа

$$u_{i+1} = u_i + u'_i(x)h + \frac{1}{2}u''(x_i + \theta_i h)h^2. \quad (47)$$

Подставляя разложение (47) в формулу (34) для погрешности аппроксимации уравнения, получим

$$\psi_i = \frac{1}{2}u''(x_i + \theta_i h)h. \quad (48)$$

Согласно формуле (46) функция $u''(x)$ непрерывна и ограничена

$$|u''(x)| \leq M_2, \quad x_i \leq x \leq x_{i+1}. \quad (49)$$

Это позволяет написать оценку

$$|\psi_i| \leq \frac{M_2}{2} h, \quad \|\psi\| \leq \frac{M_2}{2} h e^{Ch}. \quad (50)$$

Неравенства (50) показывают, что при $h \rightarrow 0$ погрешность аппроксимации уравнения и связанный с ней неравенством (45) погрешность решения стремятся к нулю со скоростью h . В связи с этим метод Эйлера называют методом первого порядка точности относительно h .

Задача 3.

Рассмотреть задачу Коши

$$u' = \frac{1}{2}u + x, \quad (51)$$

$$u(0) = 0. \quad (52)$$

Построить ее численное решение на отрезке $[0, 2]$ по схеме Эйлера с шагами $h = 0.25$, $h = 0.5$, $h = 1$. Сравнить результаты расчетов между собой и с аналитическим решением задачи

$$u(x) = -2(x+2) + 4e^{x^2}. \quad (53)$$

Результаты расчетов приведены в таблице 1.

Наиболее удобные разностные схемы этого семейства соответствуют двум значениям параметра α : $\alpha = \frac{1}{2}$ и $\alpha = 1$. При $\alpha = \frac{1}{2}$ рекуррентная формула (62) принимает вид

$$y_{i+1} = y_i + \frac{h}{2} \left\{ f(x_i, y_i) + f(x_i + h, y_i + hf(x_i, y_i)) \right\}. \quad (63)$$

Она определяет следующую процедуру расчета y_{i+1} . Сначала делается шаг h по схеме Эйлера и вычисляется величина

$$\bar{y}_{i+1} = y_i + f(x_i, y_i)h. \quad (64)$$

Затем находится значение функции f в точке (x_{i+1}, \bar{y}_{i+1}) , составляется полу сумма

$$\frac{f(x_i, y_i) + f(x_{i+1}, \bar{y}_{i+1})}{2}$$

и проводится окончательный расчет величиной

$$y_{i+1} = y_i + \frac{f(x_i, y_i) + f(x_{i+1}, \bar{y}_{i+1})}{2} h. \quad (65)$$

Такая схема вычислений называется «предиктор-корректор» или буквально «предсказание-исправление». Вычисление \bar{y}_{i+1} по схеме Эйлера – это грубое предсказание результата. Вторичный расчет (65), сделанный на основании первого, является уточнением результата, его коррекцией.

При $\alpha = 1$ рекуррентная формула (62) имеет вид

$$y_{i+1} = y_i + \frac{h}{2} \left\{ y_i + \frac{h}{2} f(x_i, y_i) \right\}. \quad (66)$$

Здесь схема расчета заключается в следующем. Сначала делается половинный шаг $\frac{h}{2}$ по схеме Эйлера вычисляется величина

$$y_{i+\frac{1}{2}} = y_i + \frac{h}{2} f(x_i, y_i). \quad (67)$$

Затем находится значение функции f в точке $(x_{i+\frac{1}{2}}, y_{i+\frac{1}{2}})$. Оно определяет по формуле (66) очередное значение $y_{i+\frac{1}{2}}$.

Следует заметить, что процедура расчета приближенного решения задачи Коши (27), (28) по схеме (61) по сравнению со схемой Эйлера усложняется: теперь на каждом шаге функцию $f(x, u)$ приходится считать не один, а два раза. Однако такое усложнение оказывается оправдано благодаря более высокой точности метода. К исследованию проблемы точности мы теперь и переходим.

Введем, как и в предыдущем разделе, две сеточные функции: погрешность решения \mathbf{z} (33) и погрешность аппроксимации уравнения ψ . В рассматриваемом случае она определяется формулой

$$\psi_i = \frac{u_{i+1} - u_i}{h} - \left[(1-\alpha) f(x_i, u_i) + \alpha f\left(x_i + \frac{h}{2}, u_i + \frac{h}{2} f(x_i, u_i)\right)\right]. \quad (68)$$

Выразим y_i по формуле (35) через u_i и z_i и подставим в разностное уравнение (56). В результате получим

$$\frac{z_{i+1} - z_i - u_i - u_{i-1}}{h} = (1-\alpha)f(x_i, u_i + z_i) + \alpha f\left(x_i + \frac{h}{2\alpha}, u_i + z_i + \frac{h}{2\alpha}f(x_i, u_i + z_i)\right). \quad (69)$$

Формулу (69) можно переписать в виде

$$\frac{z_{i+1} - z_i}{h} = \left[\left((1-\alpha)f(x_i, u_i + z_i) + \alpha f\left(x_i + \frac{h}{2\alpha}, u_i + z_i + \frac{h}{2\alpha}f(x_i, u_i + z_i)\right) \right) - \left((1-\alpha)f(x_i, u_i) + \alpha f\left(x_i + \frac{h}{2\alpha}, u_i + \frac{h}{2\alpha}f(x_i, u_i)\right) \right) \right] - \left[\frac{u_{i+1} - u_i}{h} - \left((1-\alpha)f(x_i, u_i) + \alpha f\left(x_i + \frac{h}{2\alpha}, u_i + \frac{h}{2\alpha}f(x_i, u_i)\right) \right) \right]. \quad (70)$$

Здесь мы перенесли член $\frac{u_{i+1} - u_i}{h}$ влево направо и в каждом из двух выражений, собранных в фигурных скобках, добавили одно и то же слагаемое. Поскольку между фигурными скобками стоит знак минус, значение правой части формулы (70) в целом при этом не меняется. Однако благодаря таким преобразованиям мы собрали во вторых фигурных скобках члены, которые дают погрешность аппроксимации дифференциального уравнения ψ (68).

Перейдем к дальнейшему исследованию соотношения (70). Рассмотрим функцию

$$F(v) = (1-\alpha)f(x_i, v) + \alpha f\left(x_i + \frac{h}{2\alpha}, v + \frac{h}{2\alpha}f(x_i, v)\right). \quad (71)$$

Выражение, стоящее в первых фигурных скобках формулы (70), можно записать как разность значений этой функции при $v = u_i + z_i$ и $v = u_i$ и преобразовать эту разность с помощью формулы конечных приращений Лагранжа

$$F(u_i + z_i) - F(u_i) = F'(u_i + \theta z_i)z_i, \quad 0 < \theta < 1, \quad (72)$$

где

$$F'(v) = (1-\alpha)\frac{\partial f}{\partial v}(x_i, v) + \alpha \frac{\partial f}{\partial v}\left(x_i + \frac{h}{2\alpha}, v + \frac{h}{2\alpha}f(x_i, v)\right) + \left(1 + \frac{h}{2\alpha}\frac{\partial^2 f}{\partial v^2}(x_i, v)\right). \quad (73)$$

Подставим полученные выражения для отдельных слагаемых в формулу (70). В результате она примет вид рекуррентной формулы

$$z_{i+1} = [1 + hF'(u_i + \theta z_i)]z_i - \psi_i h, \quad 0 \leq i \leq n-1, \quad (74)$$

которую нужно дополнить нулевым начальным условием (38). Использовать эту формулу для последовательного вычисления значений сеточной функции z нельзя: в ее правую часть входят неизвестные аргументы: u_i , z_i . Однако эту систему рекуррентных равенств можно заменить системой рекуррентных неравенств для последующей оценки z_i .

в столбцах приведены результаты расчетов по методу Адамса. Они будут обсуждаться в следующем разделе.

Таблица 2.

| x_i | Р.К. - II | Р.К. - IV | Ад. - II | Ад. - IV | $u(x_i)$ |
|-------|-----------|-----------|----------|----------|----------|
| 0.00 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 0.25 | 0.031250 | 0.032593 | 0.031250 | 0.032593 | 0.032594 |
| 0.50 | 0.133057 | 0.136099 | 0.130859 | 0.136099 | 0.136102 |
| 0.75 | 0.314197 | 0.319962 | 0.309692 | 0.319962 | 0.319966 |
| 1.00 | 0.587068 | 0.594879 | 0.578331 | 0.594879 | 0.594883 |
| 1.25 | 0.961913 | 0.972975 | 0.948662 | 0.972975 | 0.972984 |
| 1.50 | 1.452948 | 1.467988 | 1.434141 | 1.467972 | 1.468000 |
| 1.75 | 2.075605 | 2.095486 | 2.050001 | 2.095159 | 2.095501 |
| 2.00 | 2.847365 | 2.873107 | 2.813492 | 2.872644 | 2.873127 |

Сравнение результатов второго столбца таблицы 1, рассчитанных по методу Эйлера с шагом $h=0.25$, с результатами второго и третьего столбца таблицы 2 показывает как уменьшается погрешность при фиксированном шаге h на мере перехода к более точным методам. Так метод Рунге-Кутта четвертого порядка, несмотря на достаточно крупный шаг,ает погрешность $\|\psi\|=0.00002$. Это на много лучше, чем при расчете по схеме Эйлера с шагом $h=0.01$ (см. четвертый столбец в таблице 1). В то же время при расчете по схеме Эйлера было сделано двести шагов с однократным вычислением функции $f(x, y)$ на каждом шаге, а при расчете по схеме Рунге-Кутта – восемь шагов с четырехкратным вычислением функции $f(x, y)$ на каждом шаге. Таким образом, более сложный, но и более совершенный метод позволяет при меньшем числе вычислений получить более точный результат.

В заключение сделаем следующее замечание. Априорные оценки погрешности по схеме Эйлера (50) или Рунге-Кутты (83) представляют теоретический интерес. Они определяют скорость, с которой погрешность стремится к нулю при $h \rightarrow 0$. Однако на практике оценки подобного типа неэффективны, поскольку содержат производные искомого решения $y(x)$. Однако точность численного решения задачи устанавливается с помощью апостериорных оценок, основанных на сравнении результатов расчетов с шагом h и $h/2$. Процедура их вывода и применения была описана в предыдущей главе в связи с задачей численного интегрирования.

2.4. Метод Адамса.

Адамс – английский астроном и математик XIX века, который много занимался небесной механикой. При изучении траекторий планет ему постоянно приходилось численно интегрировать уравнения их движения. Желая минимизировать объем вычислений, Адамс разработал один из наиболее экономичных методов численного решения дифференциальных уравнений, к обсуждению которого мы теперь переходим.

Пусть $y(x)$ – решение дифференциального уравнения (27). Для производной этой функции имеет место равенство

$$|y''(x)| \leq M_3, \quad x_0 \leq x \leq x_0 + h. \quad (99)$$

Мы видим, что разностное уравнение метода Адамса, соответствующее случаю $m=1$, аппроксимирует дифференциальное уравнение (27) со вторым порядком точности относительно h . Как и в случае метода Рунге-Кутты, это обеспечивает второй порядок точности для погрешности решения $\|\psi\|$ при предположении, что значение y_i , которое рассчитывается нестандартно, вычислено со вторым порядком точности.

Процесс построения более точных схем можно продолжить за счет увеличения m . При $m=2$ получается схема третьего порядка точности, при $m=3$ – четвертого и т.д. Схема четвертого порядка, как и метод Рунге-Кутты, является наиболее удобительной, поэтому мы коротко остановимся на ее выводе и обсуждении.

Если написать интерполяционный полином третьей степени $P_3(x)$ (89) на сектке из четырех точек $x_i, x_{i-1}, x_{i-2}, x_{i-3}$ и провести интегрирование (92), то рекуррентная формула (91) примет вид

$$y_{i+1} = y_i + h \left[\frac{55}{24}f(x_i, y_i) - \frac{59}{24}f(x_{i-1}, y_{i-1}) + \frac{37}{24}f(x_{i-2}, y_{i-2}) - \frac{9}{24}f(x_{i-3}, y_{i-3}) \right]. \quad (100)$$

Приведем еще одно формулирование этой формулы через так называемые конечные разности

$$y_{i+1} = y_i + h f'_i + \frac{1}{2}h^2 A f_i + \frac{5}{12}h^3 B f'_i + \frac{3}{8}h^4 C f''_i, \quad (101)$$

где

$$f'_i = f(x_i, y_i) - f(x_{i-1}, y_{i-1}), \quad (102)$$

$$A f_i = \frac{1}{h} \{ f(x_i, y_i) - 2f(x_{i-1}, y_{i-1}) + f(x_{i-2}, y_{i-2}) \}, \quad (102)$$

$$B f'_i = \frac{1}{h^2} \{ f(x_i, y_i) - 3f(x_{i-1}, y_{i-1}) + 3f(x_{i-2}, y_{i-2}) - f(x_{i-3}, y_{i-3}) \}. \quad (102)$$

Первая, вторая и третья разности (102) приближенно соответствуют первой, второй и третьей производной функции $F(x) = f(x, y(x))$. Эквивалентность формул (100) и (101) легко проверить непосредственно. Формула (101) иногда более удобна для организации вычислительного процесса и контроля точности.

Особенность метода Адамса проявляется в формуле (100) еще сильнее, чем в формуле (93). Здесь для расчета опережающего значения y_{i+1} нужно знать значения y в четырех предыдущих точках – $y_i, y_{i-1}, y_{i-2}, y_{i-3}$. Таким образом, формула (100) начинает работать только с четвертым шагом. Вычислить по ней y_1, y_2, y_3 нельзя. Эти значения решения разностной задачи приходится рассчитывать другим методом, например, методом Рунге-Кутты.

Перейдем к обсуждению точности схемы (100). Если функция $f(x, y)$ имеет непрерывные четвертые производные по своим аргументам в интересующей нас области их изменения, так что решение задачи $u(x)$ пять раз непрерывно

Предположим, как и при исследовании метода Эйлера, что частная производная $\frac{\partial f}{\partial u}(x, u)$ в интересующей нас области изменения ее аргументов ограничена (40). Тогда с учетом формулы (73) для производной $F'(v)$ получим

$$|1 + hF'(u_i + \theta z_i)| \leq 1 + Ch + \frac{1}{2}Ch^2 < e^{\lambda h} = q, \quad q > 1. \quad (75)$$

С учетом этого рекуррентные равенства (74) можно заменить рекуррентными неравенствами

$$|z_{i+1}| \leq |z_i| + q|z_i| + \|\psi\| h, \quad (76)$$

которые полностью совпадают с неравенствами (42) предыдущего раздела. Мы уже знаем, что из них следует оценка нормы погрешности решения через норму погрешности аппроксимации уравнения

$$|\psi_i| \leq c \|\psi\|. \quad (77)$$

Теперь нужно оценить норму погрешности аппроксимации уравнения (68). Предположим, что функция $f(x, u)$ имеет в интересующей нас области изменения своих аргументов непрерывные вторые производные и, следовательно, решение дифференциального уравнения $u(x)$ трижды непрерывно дифференцируемо. Это позволяет написать следующие разложения Тейлора

$$u_{i+1} = u(x_i + h) = u_i + u'(x_i)h + \frac{1}{2}u''(x_i)h^2 + \frac{1}{6}u'''(x_i)h^3. \quad (78)$$

$$f\left(x_i + \frac{h}{2\alpha}, u_i + \frac{h}{2\alpha}f(x_i, u_i)\right) = f(x_i, u_i) + \frac{h}{2\alpha} \left\{ \frac{\partial f}{\partial x}(x_i, u_i) + \frac{\partial f}{\partial u}(x_i, u_i) f(x_i, u_i) \right\} + \frac{h^2}{8\alpha^2} \left\{ \frac{\partial^2 f}{\partial x^2}(x_i, u_i) + 2 \frac{\partial^2 f}{\partial x \partial u}(x_i, u_i) f(x_i, u_i) + \frac{\partial^2 f}{\partial u^2}(x_i, u_i) f^2(x_i, u_i) \right\}, \quad (79)$$

где

$$\tilde{x}_i = x_i + \theta h, \quad \tilde{u}_i = x_i + \theta \frac{h}{2\alpha}, \quad \tilde{u}_i = x_i + \theta \frac{h}{2} f(x_i, u_i), \quad 0 < \theta < 1. \quad (79)$$

Здесь последние слагаемые в обоих разложениях представляют собой остаточные члены в формуле Лагранжа, которые берутся в неизвестных нам промежуточных точках.

Полставим разложение (78), (79) в формулу (68) для погрешности аппроксимации дифференциального уравнения (27) и примем во внимание соотношения, вытекающие из этого уравнения

$$u''(x_i) = f(x_i, u_i), \quad (80)$$

$$u''(x_i) = \frac{\partial f}{\partial x}(x_i, u_i) + \frac{\partial f}{\partial u}(x_i, u_i) f(x_i, u_i). \quad (80)$$

Благодаря (80) члены первого и второго порядков относительно h сокращаются и остаются только члены второго порядка, обозначенные своим происхождением остаточным членом в разложениях (78), (79). В результате получается следующее представление для погрешности аппроксимации уравнения

$$u''(x) = f(x, u(x)) = F(x). \quad (86)$$

Интегрируя его между двумя точками сетки, получим соотношение

$$u_{i+1} = u_i + \int_{x_{i-m}}^{x_i} f(x) dx. \quad (87)$$

Мы не можем использовать это соотношение непосредственно для перехода в процесс решения задачи от i -ой точки сетки к $(i+1)$ -ой, поскольку функция $F(x)$ нам не известна. Чтобы сделать следующий шаг, нужно приблизенно заменить эту функцию на такую функцию, которую можно вычислить. Опишем, как эта проблема решается в методе Адамса.

Пусть в процессе численного решения задачи мы довели расчет до точки x_i . В результате проведенных расчетов нам оказались известными величины y_j и $f(x_j, y_j)$, $0 \leq j \leq i$. Возьмем некоторое фиксированное целое число $m \leq i$ и построим интерполяционный многочлен m -ой степени, принимающий в точках $x_j, j = i-m \dots i$ значение $f(x_j, u_j)$

$$P_m(x) = f(x_i, u_i), \quad i-m \leq j \leq i. \quad (88)$$

Его можно записать по формуле Лагранжа

$$P_m(x) = \sum_{j=i-m}^i f(x_j, y_j) Q_{m,j}(x), \quad (89)$$

где $Q_{m,j}(x)$ специальные многочлены вида

$$Q_{m,j}(x) = \frac{(x - x_{i-m}) \cdots (x - x_{j+1})(x - x_{j-1}) \cdots (x - x_0)}{(x_{i-m} - x_{i-m}) \cdots (x_{j+1} - x_{j+1})(x_{j-1} - x_{j-1}) \cdots (x_0 - x_0)}, \quad (90)$$

которые мы уже рассматривали в третьей главе.

Главная идея метода Адамса заключается в том, чтобы для расчета y_{i+1} использовать формулу типа (87), приближенное заменяя в ней функцию $F(x)$ на интерполяционный многочлен $P_m(x)$, составленный согласно (89) по результатам предыдущих вычислений. Это приводит к рекуррентному формуле

$$y_{i+1} = y_i + \int_{x_{i-m}}^{x_i} P_m(x) dx = y_i + \sum_{j=i-m}^i a_j f(x_j, y_j), \quad (91)$$

где

$$a_j = \int_{x_j}^{x_i} Q_{m,j}(x) dx. \quad (92)$$

Рассмотрим более подробно данную схему численного решения задачи Коши в простейших случаях $m=0$ и $m=1$, когда технические трудности не закрывают возможности идеи метода. При $m=0$ для аппроксимации функции $F(x)$ используется полином нулевой степени, т. е. постоянина

$$F(x) \approx P_0 = f(x_i, y_i).$$

дифференцируемо, то разностное уравнение (100) аппроксимирует дифференциальное уравнение (27) с четвертым порядком точности относительно h . Доказательство этого утверждения проводится также, как и для схемы второго порядка (93), только теперь в разложениях типа (90) нужно удалять больше членов. Четвертый порядок точности при аппроксимации уравнения обеспечивает четвертый порядок точности для погрешности решения $\|\psi\|$ при предположении, что начальные значения для метода Адамса y_1, y_2, y_3 вычислены с такой же точностью. Они рассчитываются независимо и при этом важно, чтобы начальный этап вычислительного процесса не внес такую погрешность, которая исказит все последующие результаты.

Задача 5.

Построить решение задачи Коши (51), (52) на отрезке $[0, 2]$ с шагом $h=0.25$ по схеме Адамса второго (93) и четвертого (100) порядка. Сравнить результаты расчетов между собой с результатами расчетов по схеме Рунге-Кутта и с аналитическим решением задачи.

Результаты расчетов приведены в четвертом и пятом столбцах таблицы 2. В соответствии с заданием, нужно сравнивать четвертый столбец со вторым и шестым, а пятый – с третьим и шестым. Напомним, что в шестом столбце приведено аналитическое решение (53) рассматриваемой задачи, так что сравнение с ним позволяет судить о точности приближенного решения по схеме Рунге-Кутты и схеме Адамса.

Расчет по схеме Адамса второго порядка точности начинается с y_2 , четвертого с y_4 . Значение y_1 в четвертом столбце, y_2, y_3, y_4 в пятом столбце рассчитываются по схеме Рунге-Кутты, соответствующему порядку, поэтому в шестом столбце они оказываются одинаковыми с соответствующими данными второго и третьего столбцов. Сравнение результатов проведенных расчетов двумя методами с аналитическим решением задачи показывает, что их точность примерно одинаковая.

Сравнение схем четвертого порядка точности в методе Рунге-Кутты (84) и Адамса (100) с точки зрения организации вычислительного процесса. Чтобы сделать шаг A и B , нужно выполнить функцию $f(x, y)$ четыре раза (85), а метод Адамса только один раз. В трех предыдущих точках функция $f(x, y)$ была уже вычислена на предыдущих шагах и вычислить ее снова нет необходимости. В этом заключается главное достоинство метода Адамса, которое особенно высоко ценится в документации.

Главный недостаток метода Адамса мы уже отмечали: при его применении первые шаги приходится делать с помощью другого метода, например, с помощью метода Рунге-Кутты и только после этого можно перейти на расчет по схеме Адамса. Таким образом, программа решения задачи Коши по методу Адамса должна включать в себя как элемент программы метода Рунге-Кутты для расчета начальной стадии вычислительного процесса.

С этой особенностью метода Адамса связана еще одна проблема. При численном интегрировании дифференциального уравнения часто приходится менять шаг h . В методе Рунге-Кутты это не составляет труда, поскольку каждый шаг делается

$$\psi_i = h^2 \left[\frac{1}{6} u''(\bar{x}) - \frac{1}{8\alpha} \left(\frac{\partial^2 f}{\partial x^2}(\bar{x}, \bar{u}) + \frac{\partial^2 f}{\partial u^2}(\bar{x}, \bar{u}) f^2(\bar{x}, \bar{u}) \right) \right]. \quad (81)$$

Функции, входящие в правую часть этого соотношения, по предположению непрерывны и ограничены в интересующей нас области изменения своих аргументов. Это позволяет заменить равенство (81) наследственным

$$|\psi_i| \leq M \|\psi\|^2. \quad (82)$$

где M – константа, мажорирующая выражение в фигурных скобках формулы (81). Подставляя оценку (82) в неравенство (77), получим

$$|\psi_i| \leq M \|\psi\|^2. \quad (83)$$

Таким образом, при $h \rightarrow 0$ погрешность аппроксимации уравнения и, как следствие, погрешность решения стремятся к нулю со скоростью h^2 . Это означает, что разностное уравнение (61), полученнное по схеме Рунге-Кутты, имеет второй порядок точности относительно h .

Второй порядок точности лучше, чем первый, однако практика показывает, что этой точности недостаточно. Наиболее удобным способом вычисления реальных расчетов используется схема Рунге-Кутты четвертого порядка точности

$$y_{i+1} = \frac{1}{6} (k_1 + 2k_2 + k_3 + k_4), \quad (84)$$

где

$$k_1 = f(x_i, y_i), \quad k_2 = f\left(x_i + \frac{h}{2}, y_i + \frac{h}{2} k_1\right), \quad (85)$$

$$k_3 = f\left(x_i + \frac{h}{2}, y_i + \frac{h}{2} k_2\right), \quad k$$

(108) сеточная функция $\{y_i\}$ определяется из решения системы линейных алгебраических уравнений. Такая разностная схема называется и явной.

Из записи разностных уравнений в форме (109) видно, что мы получили систему уравнений с трехдиагональной матрицей с диагональным преобразованием: диагональный элемент $(2+qh^2)$ больше суммы двух других элементов той же строки, равной 2. Системы такого типа мы уже встречали в третьей главе в связи с задачей интерполяции кубическим сплайном. Диагональное преобразование гарантирует существование и единственность решения системы, которое может быть построено методом прогонки.

Перейдем к обсуждению основного вопроса: с какой точностью сеточная функция $\{y_i\}$, полученная в результате решения задачи (107), (108), приближает решения краевой задачи (103), (104). Пусть $u(x)$ решение исходной краевой задачи. Обозначим через $u_i = u(x_i)$ его значения в узлах сетки и введем две сеточные функции: погрешность решения и погрешность аппроксимации уравнения

$$z_i = y_i - u_i, \quad 0 \leq i \leq n, \quad (110)$$

$$\psi_i = \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} - q_i u_i + f_i, \quad 1 \leq i \leq n-1. \quad (111)$$

Выразим из соотношения (110) y через u и z и подставим в разностное уравнение (107). Оставим члены, содержащие z_i , слева, а остальные члены перенесем направо. В результате получим

$$\frac{z_{i-1} - 2z_i + z_{i+1}}{h^2} - q_i z_i = -\left\{ \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} - q_i u_i + f_i \right\} = -\psi_i, \quad 1 \leq i \leq n-1. \quad (112)$$

Границные условия в дифференциальной и разностной задачах совпадают, так что значения сеточной функции z_i в граничных точках будут нулевыми

$$z_0 = z_n = 0. \quad (113)$$

Мы не можем рассчитать погрешность $\{z_i\}$, решая задачу (112), (113), поскольку в правые части уравнений входят неизвестные величины u_i и ψ_i . Однако задача (112), (113) позволяет оценить погрешность.

Пусть максимальное по модулю число z_i соответствует индексу $i=j$:

$$\|z\| = |z_j| \geq |z_i|, \quad 0 \leq i \leq n. \quad (114)$$

В граничных точках z_i обращается в ноль (113), так что индекс j не равен ни нулю, ни n . Рассмотрим уравнение (112) для этого значения индекса и запишем его в виде:

$$(2+qh^2)z_j = z_{j-1} + z_{j+1} + \psi_j h^2. \quad (115)$$

Возьмем модуль от обеих частей равенства и оценим правую часть сверху

$$(2+qh^2)|z_j| = (2+qh^2)\|z\| \leq |z_{j-1}| + |z_{j+1}| + |\psi_j| h^2 \leq 2\|z\| + 2\|\psi\| h^2$$

или

$$\|z\| \leq \frac{1}{q_0} \|\psi\|. \quad (116)$$

Здесь мы сократили одинаковые члены слева и справа, разделили обе части неравенства на множитель $q_0 h^2$ и заменили q_i в знаменателе на минимально возможное значение функции $q(x)$ на отрезке $[a,b]$, равное q_0 (105). Таким образом нам удалось оценить погрешность решения $\|\psi\|$ через погрешность аппроксимации уравнения $\|\Psi\|$.

Для оценки погрешности аппроксимации уравнения предположим, что функции $f(x)$ и $q(x)$ дважды непрерывно дифференцируемы на отрезке $[a,b]$. В этом случае уравнение (103) допускает двухкратное дифференцирование, что обеспечивает существование у решения краевой задачи (103), (104) четырех непрерывных производных и позволяет написать разложение

$$u_{i-1} = u_i - u'(x_i)h + \frac{1}{2}u''(x_i)h^2 - \frac{1}{6}u'''(x_i)h^3 + \frac{1}{24}u^{(4)}(x_i)h^4, \quad (117)$$

$$u_{i+1} = u_i + u'(x_i)h + \frac{1}{2}u''(x_i)h^2 + \frac{1}{6}u'''(x_i)h^3 + \frac{1}{24}u^{(4)}(x_i+\tilde{\theta}h)h^4.$$

Подставляя их в формулу (111), получим следующее выражение для ψ_i :

$$\psi_i = [u'(x_i) - q_i u_i + f_i] + \frac{h^2}{24} [u^{(4)}(x_i - \tilde{\theta}h) + u^{(4)}(x_i + \tilde{\theta}h)]. \quad (118)$$

Выражение в первых фигурных скобках равняется нулю в силу дифференциального уравнения (103). В результате в правой части формулы (118) остается только вторая группа членов, обозначенная своим происхождением остаточным членом в разложениях (117). Оценим ее следующим образом. Функция $u^{(4)}$ непрерывна и, следовательно, ограничена на отрезке $[a,b]$. Пусть

$$|u^{(4)}(x)| \leq M_4, \quad a \leq x \leq b, \quad (119)$$

тогда из формул (116) и (118) получаем

$$\|\psi\| \leq \frac{M_4}{12} h^2, \quad \|z\| \leq \frac{M_4}{12q_0} h^2. \quad (120)$$

Мы видим, что разностная схема (107) обеспечивает второй порядок аппроксимации уравнения и, как следствие неравенства (116), второй порядок точности для погрешности решения.

Задача 6.

Рассмотреть на отрезке $[-1,1]$ краевую задачу

$$u'' - u = -1, \quad (121)$$

$$u(-1) = u(1) = 0. \quad (122)$$

Выписать и решить соответствующую разностную задачу с шагом $h=0.5$.

Сравнить решение разностной задачи с аналитическим решением

$$u(x) = 1 - \frac{chx}{ch1}, \quad (123)$$

Система трех уравнений относительно y_1 , y_2 , y_3 с учетом нулевых граничных условий имеет вид

$$\begin{cases} -2.25y_1 + y_2 &= -0.25 \\ y_1 - 2.25y_2 + y_3 &= -0.25 \\ y_2 - 2.25y_3 &= -0.25 \end{cases} \quad (124)$$

Решение системы (124), как и решение исходной дифференциальной задачи, симметрично относительно средней точки, так что $y_i = u_i$. С учетом этой особенности система (124) сводится к системе двух уравнений с двумя неизвестными:

$$-2.25y_1 + y_2 = -0.25$$

$$2y_1 - 2.25y_2 = -0.25,$$

решение которой имеет вид

$$\begin{cases} 0.8125 \\ y_1 = y_2 = 0.265306 \\ 3.0625 \end{cases} = 0.346939.$$

В таблице 3 приведены значения x_i , соответствующие узлам сетки, решение разностной задачи y_i , аналитическое решение (123), вычисленное в узлах сетки $u_i = u(x_i)$, погрешность решения z_i (110) и погрешность аппроксимации уравнения $\|\psi\|$ (118). Согласно двум последним столбцам

$$\|z\| = 0.005007, \quad \|\psi\| = 0.015352$$

(125)

Таблица 3

| x_i | y_i | $u(x_i)$ | z_i | ψ_i |
|-------|----------|----------|-----------|-----------|
| -1.0 | 0.000000 | 0.000000 | 0.000000 | |
| -0.5 | 0.265306 | 0.269237 | -0.003931 | -0.015352 |
| 0.0 | 0.346939 | 0.351946 | -0.005007 | -0.013614 |
| 0.5 | 0.265306 | 0.269237 | -0.003931 | -0.015352 |
| 1.0 | 0.000000 | 0.000000 | 0.000000 | |

Погрешность аппроксимации уравнения Ψ определена только для внутренних точек сетки, поэтому первая и последние строки последнего столбца остались незаполненными.

Теперь обратимся к теоретической оценке погрешности решения и погрешности аппроксимации уравнения (120). В данном случае

$$u^{(4)}(x) = -\frac{chx}{ch1}, \quad \text{так что } |u^{(4)}(x)| \leq M_4 = 1.$$

В результате оценки (120) с учетом того, что $q=1$, получим

$$\|z\| \leq \frac{0.25}{12} = 0.020833, \quad \|\psi\| \leq \frac{0.25}{12} = 0.020833.$$

Это согласуется с фактическими значениями погрешности (125), подсчитанными непосредственно по известному решению краевой задачи (123).