

МЕТОДЫ ВЫЧИСЛЕНИЙ

Б. В. Фалейчик

20 сентября 2020 г.

Содержание

1. Теоретический материал	4
1. Предварительные сведения	5
1.1. Машинная арифметика	6
1.2. Обусловленность задачи	13
2. Методы решения СЛАУ	19
2.1. Метод Гаусса	20
2.2. LU-разложение	29
2.3. Разложение Холецкого	35
2.4. Метод отражений	40
2.5. Задача наименьших квадратов	47
2.6. Форматы хранения разреженных матриц	50
2.7. Итерационные методы решения СЛАУ	54
2.8. Чебышевское ускорение итерационных процессов	62
2.9. Методы крыловских подпространств	68
3. Методы решения проблемы собственных значений	75
3.1. Проблема собственных значений: общая характеристика	76
3.2. Степенной метод	77
3.3. Метод Данилевского	84
3.4. QR -алгоритм	90
3.5. Метод вращений Якоби	97
4. Методы решения нелинейных уравнений и систем	103
4.1. Решение нелинейных уравнений	104
4.2. Решение систем нелинейных уравнений	114
4.3. Анализ сходимости итерационных процессов	120
4.4. Численные методы оптимизации	126

5. Приближение функций	132
5.1. Интерполяция функций. Интерполяционный многочлен .	133
5.2. Интерполяционный многочлен в форме Ньютона	139
5.3. Остаток интерполяции	144
5.4. Сплаины	152
5.5. Приближение кривых	160
5.6. Среднеквадратичные приближения	168
5.7. Приближение поверхностей	177
6. Приближенное вычисление интегралов	184
6.1. Интерполяционные квадратурные формулы	185
6.2. Остаток квадратурных формул	189
6.3. Квадратурные формулы Ньютона — Котеса	193
6.4. Квадратурные формулы Гаусса	199
6.5. Составные квадратурные формулы	206
7. Численное решение ОДУ	215
7.1. Одношаговые методы решения задачи Коши	216
7.2. Методы Рунге — Кутты	223
7.3. Выбор шага численного интегрирования ОДУ	231
7.4. Многошаговые методы	239
7.5. Численное решение краевых задач. Метод стрельбы . . .	245
7.6. Решение краевых задач конечно-разностным методом .	249
Предметный указатель	254
Список рекомендуемой литературы	256

Часть 1.

Теоретический материал

Глава 1

Предварительные сведения

1.1. Машинная арифметика

1.2. Обусловленность задачи

1.1. Машинная арифметика

1.1.1. Числа с плавающей точкой

Для того чтобы использовать вещественные числа в машинных вычислениях, необходимо решить следующую общую проблему: каким образом сохранить произвольное $x \in \mathbb{R}$ в ограниченном количестве ячеек памяти? Существует несколько способов решения этой проблемы, и наиболее распространенным является представление x в виде числа с плавающей точкой.

Определение 1.1. Пусть $\beta \in \mathbb{N}$ — основание системы счисления, $p \in \mathbb{N}$ — число значащих разрядов, d_i — цифры. Вещественное число вида

$$\pm \underbrace{d_0, d_1 d_2 \dots d_{p-1}}_m \cdot \beta^e, \quad 0 \leq d_i < \beta, \quad (1.1)$$

называется *числом с плавающей точкой* (ЧПТ). Число $m \in \mathbb{R}$ называют *мантиссой* или *значащей частью*. Число $e \in \mathbb{Z}$ называют *показателем* или *экспонентой* (не путать с числом e).

Представление (1.1) для любого x очевидно не является единственным — оно зависит от «положения точки»:

$$0,0001234 = 0,0012340 \cdot 10^{-1} = 1,2340000 \cdot 10^{-4}.$$

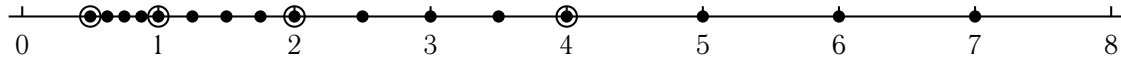
Поэтому по умолчанию используется так называемая нормализованная форма записи ЧПТ, в которой точка ставится после первой значащей цифры.

Определение 1.2. Число с плавающей точкой с ненулевым первым разрядом ($d_0 \neq 0$) называется *нормализованным*. Множество всех нормализованных ЧПТ с основанием β , p -разрядной мантиссой и $e_{\min} \leq e \leq e_{\max}$ условимся обозначать $\mathbb{F}_1(\beta, p, e_{\min}, e_{\max})$ или просто \mathbb{F}_1 .

1.1.2. Двоичные числа с плавающей точкой

Поскольку в большинстве современных компьютеров используется двоичная арифметика, рассмотрим подробно случай $\beta = 2$. Возьмем $p = 3$, $e_{\min} = -1$, $e_{\max} = 2$. Все соответствующие положительные нормализованные ЧПТ приведены в таблице, они же изображены на рисунке. Степени двойки, соответствующие единичному значению мантиссы, обведены окружностями.

$m \setminus e$	-1	0	1	10 (2)
1,00	0,1 (0,5)	1 (1)	10 (2)	100 (4)
1,01	0,101 (0,625)	1,01 (1,25)	10,1 (2,5)	101 (5)
1,10	0,110 (0,75)	1,1 (1,5)	11 (3)	110 (6)
1,11	0,111 (0,875)	1,11 (1,75)	11,1 (3,5)	111 (7)



Приведенные данные наглядно демонстрируют следующие важные свойства, общие для всех нормализованных ЧПТ:

- 1) отсутствие нуля: $0 \notin \mathbb{F}_1$;
- 2) ЧПТ распределены на числовой прямой *неравномерно*;
- 3) чем больше модуль $\xi \in \mathbb{F}_1$, тем больше и расстояние между ξ и соседними элементами \mathbb{F}_1 ;
- 4) между нулем и минимальным положительным $\xi_{\min} \in \mathbb{F}_1$ существует «зазор», ширина которого больше расстояния от ξ_{\min} до следующего ЧПТ в 2^{p-1} раз.

Двоичные ЧПТ выгодно отличаются от остальных тем, что в их нормализованной записи первый разряд d_0 всегда равен 1, поэтому его в памяти можно не хранить. *Таким образом, для хранения p -разрядной двоичной мантиссы достаточно $(p - 1)$ битов.*

1.1.3. Способы округления

Если представление x в системе счисления по основанию β содержит больше p значащих цифр, мы не можем точно записать его в виде (1.1). В этом случае можно лишь приблизить x каким-то числом с плавающей точкой, которое в дальнейшем будем обозначать $\text{fl}(x)$.

Определение 1.3. *Правилом округления* для данного множества чисел с плавающей точкой $\mathbb{F} \subset \mathbb{R}$ будем называть отображение

$$\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$$

такое, что $\text{fl}(x) = x$, если $x \in \mathbb{F}$, и $\text{fl}(x) \approx x$ в противном случае.

Рассмотрим несколько способов задания fl , считая $\beta = 10$, $p = 3$.

- 1) Отбрасывание «лишних» знаков ($\text{fl} = \text{fl}_d$): $\text{fl}_d(0,12345) = 1,23 \cdot 10^{-1}$.
- 2) Округление вверх («школьное округление», $\text{fl} = \text{fl}_u$). $\text{fl}_u(543,21) = 5,43 \cdot 10^2$, $\text{fl}_u(5678) = 5,68 \cdot 10^3$. В случае, когда запись x заканчивается на 5, его округляют до большего ЧПТ (вверх): $\text{fl}_u(23,45) = 2,35 \cdot 10^1$.

3) Округление до четного ($\text{fl} = \text{fl}_e$). Этот способ отличается от предыдущего только трактовкой «спорного» случая, когда x находится ровно между двумя ЧПТ \underline{x} и \bar{x} . Оба эти приближения на самом деле равноправны, поэтому вместо того, чтобы всегда выбирать \bar{x} , с вероятностью 50 % выбирается $\text{fl}_e(x) = \underline{x}$ либо $\text{fl}_e(x) = \bar{x}$. Реализовать это можно, всегда выбирая из \underline{x} и \bar{x} то значение, мантисса которого заканчивается на четную цифру. Таким образом, получим: $\text{fl}_e(23,45) = 2,34 \cdot 10^1$, но $\text{fl}_e(23,55) = 2,36 \cdot 10^1$.

Возникает вопрос: какой из описанных способов округления лучше? Ответ на него дает следующая теорема.

Теорема 1.1. Пусть x и y — два числа с плавающей точкой. Рассмотрим последовательность $\{x_i\}$, определенную по правилу

$$x_0 = x, \quad x_{i+1} = \text{fl}(\text{fl}(x_i + y) - y).$$

Если $\text{fl} = \text{fl}_e$, то либо $x_i = x \ \forall i \geq 0$, либо $x_i = x_1 \ \forall i \geq 1$.

Заметим, что по стандарту IEEE 754 в современных ЭВМ используется $\text{fl} = \text{fl}_e$.

1.1.4. Расширение множества чисел с плавающей точкой

Для практического использования машинной арифметики нам недостаточно множества нормализованных ЧПТ \mathbb{F}_1 . Как минимум к этому множеству нужно добавить нуль (свойство 1). Кроме этого, современные машинные арифметики включают специальные значения для обозначения бесконечностей, результатов некорректных операций и т. п.

Денормализованные числа

Наличие «зазора» между нулем и минимальным положительным нормализованным ЧПТ (свойство 4) может привести к серьезным проблемам на практике. Возьмем, например, ЧПТ $x = 0,75$ и $y = 0,625$ из рассмотренного выше модельного множества $\mathbb{F}_1(2, 3, -1, 2)$. Поскольку $x - y = 0,125$, при любом разумном способе округления мы имеем $\text{fl}(x - y) = 0$, т. е., например, выполнение обычного кода

```
if (x != y) then z = 1/(x - y)
```

в нашем случае приведет к плачевному результату.

Эта проблема в современных машинных арифметиках решается дополнением множества нормализованных ЧПТ так называемыми денормализованными числами.

Определение 1.4. Вещественные числа вида

$$0, d_1 d_2 \dots d_{p-1} \cdot \beta^{e_{\min}},$$

где d_i — произвольные цифры (по основанию β), называются *денормализованными* числами с плавающей точкой (ДЧПТ). Множество всех ДЧПТ с параметрами β, p, e_{\min} будем обозначать $\mathbb{F}_0(\beta, p, e_{\min})$ либо кратко \mathbb{F}_0 .

Введение денормализации сразу решает две проблемы: теперь выполняется свойство

$$x = y \Leftrightarrow \text{fl}(x - y) = 0 \quad \text{для любых ЧПТ } x, y,$$

а также добавляется нуль ко множеству машинных чисел.

Заметим, что для хранения денормализованных ЧПТ необходимо одно дополнительное значение для экспоненты (как правило, это $e_{\min} - 1$).

Специальные величины

Стандарт IEEE 754, которому соответствуют практически все современные ЭВМ, предусматривает наличие специальных значений для машинных чисел, которым соответствуют не ЧПТ, а другие объекты. Простейшие объекты такого типа — это $+\infty$ и $-\infty$ (присутствуют также $+0$ и -0). Результаты вычислений с бесконечностями являются вполне определенными: например, если x — положительное число, то по стандарту $x / \pm \infty = \pm 0$, $x / \pm 0 = \pm \infty$ и т. д. Кроме этого, стандартом определяются так называемые «нечисла» (NaN, от «not a number»), которые обозначают результаты некорректных арифметических операций, таких как, например, извлечение корня из отрицательного числа.

В дальнейшем под *машинными числами* будем понимать элементы множества

$$\mathbb{F}_0 \cup \mathbb{F}_1,$$

а под термином *арифметика с плавающей точкой* (АПТ) — множество машинных чисел в совокупности с правилом округления fl .

При вычислениях в АПТ будем считать, что результаты операций сложения, вычитания, умножения и деления являются *точно округляемыми*. Это означает, что результат указанных операций всегда вычисляется точно, после чего округляется до ЧПТ по правилу fl .

1.1.5. Машинный эпсилон

Параметры $\beta, p, e_{\min}, e_{\max}$ и fl полностью определяют свойства АПТ, однако их знание не дает прямой информации о том, насколько хороша или плоха соответствующая арифметика. С практической точки зрения пользователю нужны критерии, по которым можно определить качество арифметики. Основным показателем качества будем считать точность, с которой арифметика приближает вещественные числа.

Определение 1.5. *Абсолютной погрешностью округления* для числа $x \in \mathbb{R}$ в данной арифметике с плавающей точкой называется число

$$\Delta(x) = |x - \text{fl}(x)|, \quad (1.2)$$

а *относительной погрешностью округления* — число

$$\delta(x) = \frac{|x - \text{fl}(x)|}{|x|} = \frac{\Delta(x)}{|x|}. \quad (1.3)$$

Иногда относительную погрешность измеряют в процентах, умножая ее на 100. Важно понимать, что при работе с машинной арифметикой уместнее всего оперировать относительными погрешностями, так как чем больше модуль x , тем большее значение может иметь $\Delta(x)$ (свойство 3), в то время как максимальная относительная погрешность $\delta(x)$ не зависит от величины $|x|$.

Определение 1.6. *Машинным эпсилоном* (ε_M) для арифметики с плавающей точкой называется наименьшее положительное число ε , удовлетворяющее условию

$$\text{fl}(1 + \varepsilon) > 1.$$

Заметим, что значение машинного эпсилон зависит от правила округления и от количества бит в мантиссе.

▷₁ Выразите ε_M в зависимости от указанных параметров.

Определение, как это часто бывает, не указывает явно на истинный смысл величины ε_M , который выражается в следующей теореме.

Теорема 1.2. *Для всех вещественных x таких, что $\xi_{\min} \leq |x| \leq \xi_{\max}$, где ξ_{\min} и ξ_{\max} — минимальное и максимальное положительное нормализованное ЧПТ соответственно, справедлива оценка*

$$\delta(x) \leq \varepsilon_M.$$

Другими словами, величина машинного эпсилон ограничивает относительную погрешность округления для всех чисел в диапазоне нормализованных ЧПТ. Чем меньше величина ε_M , тем точнее вещественные числа приближаются в машинной арифметике.

▷₂ Докажите теорему 1.2. Для этого можно использовать следующий способ: оцените максимальную относительную погрешность $\Delta(x)$ для $x \in (1, 2)$, затем для $x \in (2, 4)$ и обобщите для $x \in (2^k, 2^{k+1})$, $k \in \mathbb{Z}$. После этого останется применить определение относительной погрешности $\delta(x)$.

▷₃ Покажите, что утверждение теоремы 1.2 не выполняется в диапазоне денормализованных ЧПТ.

1.1.6. Стандарт IEEE 754

Международный стандарт *IEEE 754 floating point standard* определяет правила организации машинной арифметики с плавающей точкой. В настоящее время ему соответствует большинство вычислительных машин. В частности, наиболее распространенный тип данных, известный как *double precision floating point* (тип `double` в C/C++), по стандарту имеет следующие параметры:

β	p	e_{\min}	e_{\max}	fl
2	53	-1022	1023	fl_e

1.1.7. Проблемы машинных вычислений

Потеря значимости. Эта проблема возникает при вычитании двух близких чисел, которые не являются точно представимыми в виде ЧПТ.

Покажем это на примере модельной арифметики с $\mathbb{F}_1(2, 3, -1, 2)$ и $\text{fl} = \text{fl}_u$: пусть $x = 4,51$, $y = 4,49$. Имеем $\text{fl}(x) = 5$, $\text{fl}(y) = 4$, и $\text{fl}(\text{fl}(x) - \text{fl}(y)) = 1$, тогда как $x - y = 0,02$. Таким образом мы получили относительную погрешность вычисления, равную 5000 %, несмотря на то, что относительная погрешность округления для x и y составляет менее 12,5 %. Отметим, что сложение этих двух чисел выполняется в данной арифметике точно.

Неассоциативность арифметических операций. Работая с машинными числами, всегда следует помнить о том, что порядок операций существенно влияет на результат. Простейший случай — нарушение привычного свойства ассоциативности: если a , b и c — машинные числа, $a \circ$ — бинарная операция, то в общем случае $\text{fl}(\text{fl}(a \circ b) \circ c) \neq \text{fl}(a \circ \text{fl}(b \circ c))$.

Итак, при использовании машинной арифметики всегда следует помнить о том, что как только в память ЭВМ записывается число x , оно автоматически превращается в число $\tilde{x} = \text{fl}(x)$, которое *почти всегда* не будет равно x . Кроме того, чем больше модуль x , тем больше может быть разница между x и \tilde{x} (абсолютная погрешность $\Delta(x)$). Относительная же погрешность округления согласно теореме 1.2 почти всегда ограничена величиной ϵ_M .

1.2. Обусловленность задачи

1.2.1. Некорректные и плохо обусловленные задачи

Постоянное присутствие ошибок округления при работе с машинной арифметикой предъявляет особые требования к вычислительным алгоритмам и требует дополнительного анализа решаемой задачи. Так как практически все числа представляются в ЭВМ с погрешностью, необходимо знать, насколько решение чувствительно к изменениям параметров задачи.

Определение 1.7. Задача называется *корректно поставленной* или просто *корректной*, если ее решение существует, единственно и непрерывно зависит от начальных данных. Если нарушено хотя бы одно из этих условий, задачу называют *некорректной*.

Решение некорректных задач на ЭВМ — весьма серьезная проблема. Если, например, нарушено условие непрерывной зависимости от параметров, то наличие малейшей погрешности в начальных данных (которая практически неминуемо произойдет, как только вы запишете эти данные в память), может кардинальным образом исказить решение.

Существует еще один класс задач, формально являющихся корректными, но решения которых, тем не менее, тоже весьма плохо ведут себя при наличии погрешностей в начальных данных — это так называемые плохо обусловленные задачи. В общих чертах плохо обусловленной называется задача, которая при маленькой *относительной* погрешности в начальных данных дает большую *относительную* погрешность в решении. Упор на относительную погрешность делается потому, что, как мы знаем из п. 1.1, при округлении вещественного числа x до машинного $\text{fl}(x)$ абсолютная погрешность $\Delta(x)$ зависит от величины $|x|$, в то время как относительная погрешность $\delta(x)$ постоянна для данной машинной арифметики.

Пример 1.1 (вычисление значения многочлена). Исследовать обусловленность задачи вычисления значения многочлена относительно погрешности, вносимой в один из его коэффициентов.

Решение. Пусть $P(a_0, \dots, a_n, x) = \sum_{i=0}^n a_i x^i$. Рассмотрим задачу вычисления значения данного многочлена, считая фиксированными x и все коэффициенты a_i кроме какого-то одного — a_k , который и будем считать параметром. Пусть вместо точного значения a_k используется «возмущенное» значение \tilde{a}_k , т. е. имеем относительную погрешность в

начальных данных, равную

$$\frac{|a_k - \tilde{a}_k|}{|a_k|}.$$

Вычислим относительную погрешность решения и выразим ее через относительную погрешность начальных данных:

$$\begin{aligned} & \frac{|P(a_0, \dots, a_n, x) - P(a_0, \dots, \tilde{a}_k, \dots, a_n, x)|}{|P(a_0, \dots, a_n, x)|} = \\ &= \frac{|(a_k - \tilde{a}_k)x^k|}{|\sum_{i=0}^n a_i x^i|} = \frac{|a_k x^k|}{|\sum_{i=0}^n a_i x^i|} \frac{|a_k - \tilde{a}_k|}{|a_k|}. \end{aligned}$$

Коэффициент

$$\varkappa = \frac{|a_k x^k|}{|\sum_{i=0}^n a_i x^i|}$$

называется *числом обусловленности*. Проблемы при решении рассматриваемой задачи могут возникать, когда это число велико: например, когда x близко к одному из корней многочлена P , или $x \gg 1$ и k достаточно велико.

Численный пример:

$$p(x) = (x - 2)^{10} = x^{10} - 20x^9 + 180x^8 - 960x^7 + \dots - 5120x + 1024.$$

Пусть $x = 3$. Тогда $P(a_0, \dots, a_{10}, x) = 1$. Изменим коэффициент $a_9 = -20$ на 0,01 (что составляет 0,05 %): $\tilde{a}_9 = -19,99$. Тогда получим

$$P(a_0, \dots, \tilde{a}_9, a_{10}, x) = 197,83,$$

что на 19683 % больше точного значения. □

Таким образом, *число обусловленности задачи показывает, во сколько раз относительная погрешность возмущенного решения может превосходить относительную погрешность соответствующих начальных данных*.

Задача называется *плохо обусловленной*, если ее число обусловленности велико.

Замечание 1.1. Естественно, понятие «большое число обусловленности» относительно. Судить о величине обусловленности можно лишь в контексте той машинной арифметики, которая используется для вычислений, а еще точнее — от величины машинного эпсилон, так как эта величина ограничивает относительную погрешность округления.

▷₁ Пусть известно значение ε_M . Укажите при каких значениях числа обусловленности задачу можно считать плохо обусловленной в такой арифметике.

1.2.2. Векторные нормы

В дальнейшем как параметры, так и решения рассматриваемых задач будут векторами пространства \mathbb{R}^n . Для исследования обусловленности задач нужно измерять «величины» этих векторов, для чего используются *векторные нормы*. Мы будем активно пользоваться двумя векторными нормами: *максимум-нормой*

$$\|x\|_{\infty} = \max_{1 \leq i \leq n} |x_i| \quad (1.4)$$

и *евклидовой нормой*

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}. \quad (1.5)$$

Обе эти нормы являются частными случаями *p-нормы*, определяемой формулой

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

1.2.3. Матричные нормы

Рассмотрим задачу решения системы линейных алгебраических уравнений (СЛАУ) вида

$$Ax = b, \quad (1.6)$$

где A — невырожденная квадратная матрица размерности n ; $x, b \in \mathbb{R}^n$. Прежде чем приступить к алгоритмам численного решения этой задачи, исследуем ее обусловленность.

Как и в случае векторных параметров, нам нужно будет как-то измерять «величину» матрицы A . Делать это мы будем с использованием операторных матричных норм. При работе с матрицами (по крайней мере в контексте линейной алгебры) всегда важно помнить, что любая матрица определяет *линейный оператор*, т. е. отображение $A : \mathbb{R}^m \rightarrow \mathbb{R}^n$, которое обладает свойством линейности

$$A(\alpha x + \beta y) = \alpha Ax + \beta Ay, \quad \forall x, y \in \mathbb{R}^n, \alpha, \beta \in \mathbb{R}.$$

Важность такой точки зрения следует хотя бы из того, что так называемое правило умножения матриц, которое обычно вводится без обоснования,

есть не что иное, как алгоритм вычисления композиции линейных операторов. В таком контексте вопрос об определении «величины» матрицы сводится к определению нормы соответствующего линейного оператора.

Определение 1.8. Нормой линейного оператора (матрицы) A называют число

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|.$$

Норма оператора полностью определяется векторной нормой, т. е. каждая векторная норма порождает (*индуцирует*) соответствующую ей операторную матричную форму (в этом случае говорят также, что матричная норма *подчинена* векторной).

В дальнейшем мы без оговорок будем предполагать, что используемая матричная норма подчинена векторной.

Норма оператора равна максимальному «коэффициенту растяжения»: она показывает, во сколько раз под его действием может увеличиться норма вектора. Поэтому *по определению* для любой операторной нормы имеем важное свойство

$$\|Ax\| \leq \|A\| \|x\|. \quad (1.7)$$

Напомним, как вычисляются матричные нормы, индуцированные векторными нормами $\|\cdot\|_\infty$ и $\|\cdot\|_2$.

- Векторной максимум-нормой $\|\cdot\|_\infty$ индуцируется матричная норма, вычисляемая по правилу

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|. \quad (1.8)$$

Данную норму будем называть (строчной) *максимум-нормой*, или *кубической* матричной нормой.

- Евклидовой векторной нормой $\|\cdot\|_2$ индуцируется матричная норма, вычисляемая по правилу

$$\|A\|_2 = \max_{1 \leq i \leq n} \sqrt{\lambda_i}, \quad (1.9)$$

где $\{\lambda_i\}_{i=1}^n$ — собственные значения матрицы A^*A . Эту норму называют *спектральной* матричной нормой.

1.2.4. Число обусловленности матрицы

Рассмотрим СЛАУ (1.6). Решение этой системы эквивалентно вычислению

$$x = A^{-1}b.$$

Исследуем обусловленность этой задачи, считая параметром вектор правой части b . Действуем по схеме, аналогичной примеру 1.2.1. Относительная погрешность начальных данных имеет вид

$$\delta_b = \frac{\|b - \tilde{b}\|}{\|b\|},$$

а возмущенного решения —

$$\delta_x = \frac{\|x - \tilde{x}\|}{\|x\|},$$

где $x = A^{-1}b$, $\tilde{x} = A^{-1}\tilde{b}$. Наша задача — установить связь между этими двумя погрешностями. Имеем

$$\delta_x = \frac{\|A^{-1}(b - \tilde{b})\|}{\|A^{-1}b\|} \leq \frac{\|A^{-1}\| \|b - \tilde{b}\|}{\|A^{-1}b\|} \leq \|A^{-1}\| \|A\| \frac{\|b - \tilde{b}\|}{\|b\|}.$$

Для получения последней оценки мы использовали тот факт, что

$$\|b\| = \|AA^{-1}b\| \leq \|A\| \|A^{-1}b\| \Rightarrow \|A^{-1}b\| \geq \|A\|^{-1} \|b\|.$$

В итоге мы получим

$$\delta_x \leq \kappa(A) \delta_b, \tag{1.10}$$

где $\kappa(A) = \|A^{-1}\| \|A\|$.

Исследование обусловленности относительно погрешностей в матрице A требует более тонкого подхода, но приводит к аналогичному результату. Таким образом, из (1.10) вытекает следующее определение.

Определение 1.9. Числом обусловленности невырожденной матрицы A называется число

$$\kappa(A) = \|A\| \|A^{-1}\|. \tag{1.11}$$

Если матрица A вырождена, ее число обусловленности полагается равным бесконечности.

Замечание 1.2. Несмотря на то, что это определение ассоциировано с матрицей A , необходимо четко понимать, что речь идет об обусловленности задачи решения СЛАУ.

Замечание 1.3. Число обусловленности по определению зависит от нормы. В случаях, когда это необходимо, мы будем употреблять говорящие обозначения $\kappa_2(A)$ и $\kappa_\infty(A)$.

Число обусловленности матрицы обладает следующими свойствами:

- 1) $\kappa(A) \geq 1 : 1 = \|A^{-1}A\| \leq \|A\| \|A^{-1}\|$;
- 2) $\kappa(AB) \leq \kappa(A)\kappa(B) : \|AB\| \|(AB)^{-1}\| \leq \|A\| \|A^{-1}\| \|B\| \|B^{-1}\|$;
- 3) если $A = A^*$, то $\kappa_2(A) = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}$, где λ_{\min} и λ_{\max} — минимальное и максимальное по модулю собственные значения матрицы A соответственно.

Замечание 1.4. Отметим, что в общем случае отсутствует прямая связь между величинами собственных значений и числом обусловленности. Например, собственные значения матрицы $A = \begin{bmatrix} 1 & \alpha \\ 0 & 1 \end{bmatrix}$ равны 1, в то время как $\kappa_2(A) \rightarrow \infty$ при $|\alpha| \rightarrow \infty$.

Замечание 1.5. Укажем на одно часто встречающееся заблуждение. Так как $\kappa(A)$ является своеобразным индикатором близости матрицы A к вырожденной, может возникнуть впечатление, что чем меньше определитель, тем больше число обусловленности. На самом же деле такой связи нет: достаточно заметить, что у A и A^{-1} взаимно обратные определители, но одинаковые числа обусловленности.

В заключение данного раздела повторим основные тезисы.

- Число обусловленности задачи показывает, во сколько раз *относительная* погрешность решения может превышать *относительную* погрешность в начальных данных.
- Величина числа обусловленности, при которой задачу можно считать плохо обусловленной, зависит от параметров используемой машинной арифметики.
- При решении плохо обусловленной задачи обычными методами *нельзя рассчитывать на получение адекватного решения.*

Глава 2

Методы решения СЛАУ

- 2.1. Метод Гаусса
- 2.2. LU-разложение
- 2.3. Разложение Холецкого
- 2.4. Метод отражений
- 2.5. Задача наименьших квадратов
- 2.6. Форматы хранения разреженных матриц
- 2.7. Итерационные методы решения СЛАУ
- 2.8. Чебышевское ускорение итерационных процессов
- 2.9. Методы крыловских подпространств

2.1. Метод Гаусса

2.1.1. Введение

Приступим теперь к изучению методов решения систем линейных алгебраических уравнений (СЛАУ)

[illegible]

которые будем записывать в матричном виде

$$Ax = b, \quad (2.1)$$

где $A = (a_{ij})$ — квадратная матрица размерности n ; $x = (x_1, \dots, x_n)^T$ — искомый вектор решения; $b = (b_1, \dots, b_n)^T$ — вектор правой части. Везде в дальнейшем будем предполагать, что задача корректно поставлена, т. е. $\det A \neq 0$.

Начнем с методов, которые позволяют найти точное решение системы за конечное число операций (при условии, что все вычисления выполняются точно). Такие методы называются *прямыми* или *точными*.

Все прямые методы решения СЛАУ базируются на простой идее: исходную систему преобразуют к эквивалентной системе (т. е. к системе с тем же решением)

$$Vx = g, \quad (2.2)$$

решение которой находится легко. Например, в наиболее простом случае $V = I$ сразу получаем $x = g$. Но чаще всего исходную систему (2.1) приводят к системе с треугольной матрицей. Так, в случае верхнетреугольной матрицы V система (2.2) имеет вид

[illegible]

Решается такая система путем последовательного выражения x_i через уже известные значения, начиная с x_n :

$$x_i = \frac{1}{v_{ii}} \left(g_i - \sum_{j=i+1}^n v_{ij} x_j \right), \quad i = n, n-1, \dots, 2, 1. \quad (2.3)$$

Вычислительный процесс (2.3) называется *обратной подстановкой* (*обратным ходом*). Заметим, что при $i = n$ верхний предел суммирования будет меньше нижнего, а в таких случаях значение суммы полагается равным нулю. Это соглашение будет нами использоваться в дальнейшем по умолчанию.

Рассмотрим теперь, каким образом осуществляется переход от исходной системы (2.1) к эквивалентной (2.2). Чаще всего для этого последовательно применяют к обеим частям (2.1) некоторые линейные преобразования T_k (умножают обе части системы $Ax = b$ на матрицы T_k):

$$\underbrace{T_N \dots T_2 T_1 A}_V x = \underbrace{T_N \dots T_2 T_1 b}_g.$$

Такой процесс называют *прямым ходом*.

Если в качестве T_k использовать элементарные преобразования, то получим самый известный прямой метод — *метод Гаусса*.

2.1.2. Базовый метод Гаусса

Обозначим \underline{a}_i i -ю строчку матрицы A . Тогда прямой ход метода Гаусса, заключающийся в приведении матрицы системы к верхнетреугольному виду с помощью элементарных преобразований, можно записать в виде следующего алгоритма.

Базовый алгоритм метода Гаусса

```

1: for  $k = \overline{1, n-1}$  do
2:   for  $i = \overline{k+1, n}$  do
3:      $l \leftarrow a_{ik}/a_{kk}$ 
4:      $\underline{a}_i \leftarrow \underline{a}_i - l \underline{a}_k$ 
5:      $b_i \leftarrow b_i - l b_k$ 
6:   end for
7: end for
```

Этап алгоритма, определяемый строками 2–4, будем называть *k*-м шагом метода Гаусса. На этом этапе с помощью элементарных преобразований обнуляются элементы *k*-го столбца, находящиеся ниже главной диагонали. Матрицу системы перед выполнением *k*-го шага будем обозначать $A^{(k)}$ ($A^{(1)} = A$). Переход от матрицы $A^{(k)}$ к $A^{(k+1)}$ можно представить в виде $A^{(k+1)} = L_k A^{(k)}$, где матрица преобразования L_k имеет следующую структуру:

Умножение произвольной матрицы M слева на матрицу L_k равносильно добавлению к i -й строке матрицы M k -й строки, умноженной на $l_i^{(k)}$ для всех i от $k + 1$ до n . Согласно алгоритму метода Гаусса (см. строку 4), имеем

где $a_{ij}^{(k)}$ — элементы матрицы $A^{(k)}$.

Пусть $[A]_k$ — матрица, составленная из k первых строк и k первых столбцов матрицы A .

Теорема 2.1. *Базовый алгоритм метода Гаусса осуществим тогда и только тогда, когда все главные угловые миноры матрицы A не равны нулю: $|[A]_k| \neq 0 \quad \forall k = \overline{1, n}$.*

Доказательство. Воспользуемся методом математической индукции. Преобразование L_1 корректно определено и первый шаг метода Гаусса выполним, если (и только если) $a_{11}^{(1)} = a_{11} = |[A]_1| \neq 0$.

Пусть выполнимо k шагов. Это означает, что существует матрица \tilde{L}_k ,

$$\tilde{L}_k = L_k L_{k-1} \dots L_1, \quad \text{и} \quad A^{(k+1)} = \tilde{L}_k A.$$

Нетрудно заметить, что матрицы \tilde{L}_k имеют блочный вид

$$\left[\begin{array}{c|c} M_k & 0 \\ \hline \boxtimes & I \end{array} \right],$$

где M_k — нижнетреугольная матрица размерности k с единицами на главной диагонали; I — единичная матрица размерности $n - k$.

Запишем равенство $\tilde{L}_k A = A^{(k+1)}$ в блочном виде:

$$\underbrace{\left[\begin{array}{c|c} M_k & 0 \\ \hline \boxtimes & I \end{array} \right]}_{\tilde{L}_k} A = \underbrace{\left[\begin{array}{c|c} U_k & \boxtimes \\ \hline 0 & \boxtimes \end{array} \right]}_{A^{(k+1)}}, \quad (2.6)$$

где U_k — верхнетреугольная матрица размерности k .

Критерием осуществимости $(k + 1)$ -го шага является условие

$$\theta_{k+1} = a_{k+1,k+1}^{(k+1)} \neq 0,$$

которое гарантирует существование L_{k+1} (см. (2.4), (2.5)).

Из (2.6) имеем

$$[\tilde{L}_k]_{k+1} [A]_{k+1} = [A^{(k+1)}]_{k+1}.$$

Так как $|[\tilde{L}_k]_{k+1}| \neq 0$ и θ_{k+1} — единственный ненулевой элемент в последней строке матрицы $[A^{(k+1)}]_{k+1}$, то

$$\theta_{k+1} \neq 0 \Leftrightarrow |[A]_{k+1}| \neq 0. \quad \square$$

Трудоёмкость метода Гаусса. Прямой ход базового алгоритма метода Гаусса требует выполнения

$$\frac{n^3}{3} + O(n^2)$$

мультипликативных операций (умножения или деления) и столько же аддитивных операций (сложения или вычитания).

▷₁ Докажите это. Оцените трудоёмкость обратного хода.

2.1.3. Связь метода Гаусса и LU -разложения

Определение 2.1. LU -разложением невырожденной матрицы A называется ее представление в виде

$$A = LU,$$

где L — нижнетреугольная матрица с единицами на главной диагонали; U — верхнетреугольная матрица.

Теорема 2.2 (связь метода Гаусса и LU -разложения). *Базовый алгоритм метода Гаусса для СЛАУ (2.1) выполним тогда и только тогда, когда существует LU -разложение матрицы A .*

Доказательство. Пусть базовый алгоритм осуществим. Тогда из формулы (2.6) при $k = n - 1$ получаем $\tilde{L}_{n-1}A = A^{(n)}$, причем по построению \tilde{L}_{n-1} — нижнетреугольная с единичной главной диагональю, а $A^{(n)}$ — верхнетреугольная. Отсюда получим $A = LU$, где $L = (\tilde{L}_{n-1})^{-1}$, $U = A^{(n)}$.

Необходимость следует из теоремы 2.1: если существует LU -разложение, то нетрудно заметить, что $|[L]_k| \neq 0$ и $|[U]_k| \neq 0$. Следовательно,

$$|[A]_k| = |[L]_k| \cdot |[U]_k| \neq 0. \quad \square$$

2.1.4. Метод Гаусса с выбором главного элемента

Для того чтобы выполнение алгоритма метода Гаусса не обрывалось при наличии нулевого главного элемента (и не только поэтому), *перед каждым шагом* метода применяется процедура, называемая *выбором главного элемента*. Суть процедуры: путем перестановки строк или столбцов матрицы $A^{(k)}$ поставить на позицию (k, k) ненулевой элемент. При этом, чтобы не «испортить» структуру матрицы, можно использовать лишь последние $n - k$ строк и $n - k$ столбцов. Существует несколько способов выбора главного элемента.

По столбцу: среди элементов $a_{ik}^{(k)}$ для i от k до n выбирается ненулевой элемент $a_{i^*k}^{(k)}$, после чего переставляются местами строки k и i^* .

По строке: среди элементов $a_{kj}^{(k)}$ для j от k до n выбирается ненулевой элемент $a_{kj^*}^{(k)}$, после чего переставляются местами столбцы k и j^* .

По матрице: среди элементов $a_{ij}^{(k)}$ для i, j от k до n выбирается ненулевой элемент $a_{i^*j^*}^{(k)}$, после чего переставляются местами строки k и i^* и столбцы k и j^* .

Рассмотрим следующие вопросы:

- 1) из каких соображений следует выбирать главный элемент;
- 2) какой способ выбора главного элемента лучше использовать?

Для ответа рассмотрим еще раз матрицу L_k (2.4). Имеем

$$\kappa_{\infty}(L_k) = (1 + \max_i |l_i^{(k)}|)^2, \quad (2.7)$$

откуда с учетом свойства 2 числа обусловленности

$$\kappa_{\infty}(A^{(n)}) = \kappa_{\infty}(\tilde{L}_{n-1}A) \leq \kappa_{\infty}(A) \prod_{k=1}^{n-1} (1 + \max_i |l_i^{(k)}|)^2.$$

▷₂ Выведите эти формулы.

Таким образом, даже если $\kappa(A)$ невелико, матрица $A^{(n)}$ может стать плохо обусловленной в случае больших значений $|l_i^{(k)}|$, т. е. *сам процесс метода Гаусса может «испортить» исходную систему*.

Для исправления ситуации мы должны минимизировать величины (2.7). С учетом (2.5) получим следующие ответы:

- 1) главный элемент должен быть максимальным по модулю среди всех рассматриваемых;
- 2) выбор главного элемента по столбцу оптимален по соотношению качество/скорость.

2.1.5. Матричные уравнения

Метод Гаусса естественным образом обобщается на случай матричных уравнений вида

$$AX = B, \quad (2.8)$$

где A , как и ранее, — квадратная матрица порядка n ; B — матрица размеров $n \times m$; X — неизвестная матрица тех же размеров, что и B . Возможно два подхода к решению таких уравнений.

- 1) Система (2.8) эквивалентна набору из m СЛАУ вида

$$Ax^{(j)} = b^{(j)}, \quad j = \overline{1, m},$$

где $x^{(j)}$ и $b^{(j)}$ — столбцы матриц X и B . Решая эти m систем, найдем искомую матрицу X .

2) Матричный метод Гаусса. Для того чтобы адаптировать построенный выше алгоритм к решению матричных уравнений, достаточно строку 5 заменить на

$$\underline{b}_i \leftarrow \underline{b}_i - l \underline{b}_k$$

(здесь \underline{b}_i — строки матрицы B), а также выполнить соответственно m алгоритмов обратного хода (для каждого из столбцов $b^{(j)}$).

2.1.6. Обращение матрицы и вычисление определителя

Обращение матрицы эквивалентно решению матричного уравнения

$$AX = I,$$

где I — единичная матрица. Для решения этого уравнения могут использоваться оба описанных выше способа. Если же известно LU -разложение матрицы A , то вычислить обратную можно следующими способами [18, п. 14.3]:

1) сначала вычислить U^{-1} , после чего решить матричное уравнение $XL = U^{-1}$;

2) найти U^{-1} и L^{-1} , затем $X = A^{-1} = U^{-1}L^{-1}$.

▷₃ Постройте алгоритм обращения треугольной матрицы.

▷₄ Постройте соответствующие алгоритмы для обоих указанных выше способов обращения матрицы.

Определитель матрицы также вычисляется с помощью метода Гаусса:

$$\tilde{L}_{n-1}A = A^{(n)} \Rightarrow 1 \cdot |A| = |A^{(n)}| = a_{11}^{(n)} a_{22}^{(n)} \dots a_{nn}^{(n)}.$$

Однако при этом нужно помнить про важный нюанс: если в ходе метода переставлялись строки и столбцы, то каждая такая операция *меняла знак определителя на противоположный*. Поэтому окончательная формула такова:

$$|A| = (-1)^p a_{11}^{(n)} a_{22}^{(n)} \dots a_{nn}^{(n)}, \quad (2.9)$$

где p — количество перестановок строк и столбцов в ходе метода.

▷₅ Как вычислить определитель, если известно LU -разложение матрицы?

2.1.7. Метод прогонки

Рассмотрим СЛАУ

$$\begin{bmatrix} d_1 & e_1 & & & & \\ c_2 & d_2 & e_2 & & & \\ & c_3 & d_3 & e_3 & & \\ & & \ddots & \ddots & \ddots & \\ & & & c_{n-1} & d_{n-1} & e_{n-1} \\ & & & & c_n & d_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_{n-1} \\ b_n \end{bmatrix}. \quad (2.10)$$

Матрицы такой структуры называются *трехдиагональными*. В приложениях достаточно часто встречаются такие системы. Особая структура этой матрицы позволяет найти решение системы методом Гаусса за $O(n)$ операций. Получаемый метод называется *методом прогонки*.

Алгоритм метода прогонки

- 1: **for** $k = \overline{2, n}$ **do** // Прямой ход
- 2: $d_k \leftarrow d_k - e_{k-1}c_k/d_{k-1}$
- 3: $b_k \leftarrow b_k - b_{k-1}c_k/d_{k-1}$
- 4: **end for**
- 5: $x_n = b_n/d_n$
- 6: **for** $k = \overline{n-1, 1}$ **do** // Обратный ход
- 7: $x_k \leftarrow (b_k - e_k x_{k+1})/d_k$
- 8: **end for**

Замечание 2.1. Данный алгоритм отличается от традиционного метода прогонки, основанного на так называемом алгоритме Томаса и приведенного в большинстве учебников (см., например, [11, с. 45]). Метод Томаса также эквивалентен методу Гаусса с той лишь разницей, что в течение прямого хода метода диагональные элементы d_k становятся единичными.

При выполнении алгоритма прогонки мы лишены возможности выбора главного элемента, так как при этом нарушилась бы трехдиагональная структура матрицы A . Следовательно, метод прогонки осуществим тогда и только тогда, когда все главные миноры матрицы отличны от нуля. Существует также более простое для проверки достаточное условие осуществимости метода прогонки. Для его доказательства нам понадобятся следующие предварительные сведения.

Определение 2.2. Если элементы матрицы A удовлетворяют усло-

виям

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}| \quad \forall i = \overline{1, n}, \quad (2.11)$$

то говорят, что такая матрица обладает свойством *диагонального преобладания* (*диагонального доминирования*). Если неравенство в (2.11) строгое, говорят о *строгом диагональном преобладании*.

Теорема 2.3. *Если матрица обладает свойством строгого диагонального преобладания, то все ее главные миноры отличны от нуля.*

Следствие 2.1. *Если матрица системы (2.10) удовлетворяет условиям*

$$|d_1| > |e_1|, \quad |d_i| > |c_i| + |e_i| \quad \forall i = \overline{2, n-1}, \quad \text{и} \quad |d_n| > |c_n|,$$

то алгоритм прогонки выполним.

▷₆ На основе метода прогонки постройте экономичный алгоритм решения трехдиагональных СЛАУ с выбором главного элемента. Такой алгоритм может быть полезен в случаях, когда матрица системы не обладает диагональным преобладанием.

2.2. LU-разложение

Зачем может понадобиться LU-разложение матрицы?

Если известно разложение $A = LU$, то решение СЛАУ $Ax = b$ может быть выполнено на порядок эффективнее, чем методом Гаусса. Для этого нужно решить две системы с треугольными матрицами:

$$Ax = b \Leftrightarrow L U x = b \Leftrightarrow \begin{cases} Ly = b, \\ Ux = y. \end{cases}$$

Сначала находится y (прямой подстановкой), затем решение исходной системы — вектор x (обратной подстановкой).

▷₁ Запишите алгоритм вычисления y и x .

Нетрудно убедиться, что трудоемкость решения СЛАУ в этом случае равна $O(n^2)$. Особенно эффективен этот подход в ситуациях, когда решается несколько систем с одной и той же матрицей A : разложение строится один раз (как мы скоро увидим, его построение практически ничем не отличается от прямого хода метода Гаусса и имеет такую же трудоемкость), зато все последующие системы решаются за $O(n^2)$ операций.

2.2.1. LU-разложение без выбора главного элемента

Пусть главные угловые миноры матрицы A отличны от нуля. Тогда осуществим прямой ход метода Гаусса, представляемый в виде

$$L_{n-1} \dots L_2 L_1 A = U. \quad (2.12)$$

Матрица U уже получена, осталось разобраться как найти L . Каждая из нижнетреугольных матриц L_k в (2.12) представляет собой единичную матрицу, к которой применены те же самые элементарные преобразования, которые применялись к матрице A при обнулении элементов под главной диагональю в k -м столбце. Другими словами, L_k это матрица с единицами на главной диагонали, а остальные ненулевые ее элементы могут находиться только в k -м столбце ниже диагонали.

Из последнего тождества имеем

$$A = LU,$$

где

$$L = (L_{n-1} \dots L_2 L_1)^{-1} = L_1^{-1} L_2^{-1} \dots L_{n-1}^{-1}. \quad (2.13)$$

Нетрудно заметить, что

L_k^{-1} может быть получена из исходной матрицы L_k путем инверсии знака у поддиагональных элементов,

или

$$L_k^{-1} = 2I - L_k.$$

Для вычисления матрицы L согласно (2.13) нужно последовательно применить к L_{n-1}^{-1} элементарные преобразования L_{n-2}^{-1} , L_{n-3}^{-1} и так далее до L_1^{-1} . Учитывая структуру матриц L_k^{-1} , легко увидеть, что

k -й столбец матрицы L в точности равен k -му столбцу матрицы L_k^{-1} .

Это правило является основой приведенного ниже алгоритма построения LU -разложения матрицы. Также из него следует, что матрица L может быть записана в виде

$$L = nI - L_1 - L_2 - \dots - L_{n-1}. \quad (2.14)$$

Учитывая все вышеизложенное, мы приходим к тому, что алгоритм построения LU разложения без выбора главного элемента имеет минимальные отличия от базового метода Гаусса:

```

for  $k = \overline{1, n-1}$  do
  for  $i = \overline{k+1, n}$  do
     $a_{ik} \neq a_{kk}$ 
     $a_{i,(k+1):n} \leftarrow a_{ik} a_{k,(k+1):n}$ 
  end for
end for

```

Здесь обозначение $j_1 : j_2$ обозначает диапазон элементов с индексами от j_1 до j_2 . В частности, $a_{i,(k+1):n}$ — это элементы i -й строки матрицы A , находящиеся в столбцах с $(k+1)$ -го по n -й. В результате выполнения алгоритма верхний треугольник матрицы A , включая диагональ, будет содержать элементы матрицы U , а нижний треугольник — внедиагональные элементы матрицы L .

2.2.2. LU -разложение с выбором главного элемента

Рассмотрим теперь прямой ход метода Гаусса с выбором главного элемента по столбцу. В этом случае вместо (2.12) будем иметь

$$L_{n-1}P_{n-1} \dots L_2P_2L_1P_1A = U,$$

где P_k — матрица перестановки, которая применяется к матрице $A^{(k)}$ перед выполнением k -го шага.

Рассмотрим подробно частный случай при $n = 4$. Матрицы L_k в этом случае будут иметь общий вид

$$L_1 = \begin{bmatrix} 1 & & & \\ \alpha_1 & 1 & & \\ \beta_1 & & 1 & \\ \gamma_1 & & & 1 \end{bmatrix}, \quad L_2 = \begin{bmatrix} 1 & & & \\ & 1 & & \\ \beta_2 & 1 & & \\ \gamma_2 & & 1 & \end{bmatrix}, \quad L_3 = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & \gamma_3 & 1 \end{bmatrix}.$$

Пусть матрицы перестановок равны

$$P_1 = \Pi_{12}, \quad P_2 = \Pi_{24}, \quad P_3 = \Pi_{34},$$

где Π_{km} — единичная матрица с переставленными k -й и m -й строками. В тождестве

$$L_3P_3L_2P_2L_1P_1A = U$$

рассмотрим отдельно произведение P_2L_1 , учитывая, что умножение на матрицу Π_{km} справа соответствует перестановке k -го и m -го столбцов:

$$P_2L_1 = \Pi_{24}L_1 = \begin{bmatrix} 1 & & & \\ \gamma_1 & & & 1 \\ \beta_1 & & 1 & \\ \alpha_1 & 1 & & \end{bmatrix} = \begin{bmatrix} 1 & & & \\ \gamma_1 & 1 & & \\ \beta_1 & & 1 & \\ \alpha_1 & & & 1 \end{bmatrix} P_2 = L'_1P_2,$$

откуда получим

$$L_3P_3L_2L'_1P_2P_1A = U.$$

Матрица L'_1 представляет собой матрицу L_1 , у которой элементы (только) первого столбца переставлены в соответствии с $P_2 = \Pi_{24}$, что также можно записать как

$$L'_1 = P_2L_1P_2,$$

так как $P_2^{-1} = P_2$. Продолжая по аналогии далее, будем иметь

$$P_3 L_2 = \Pi_{34} L_2 = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & \gamma_2 & & 1 \\ & \beta_2 & 1 & \end{bmatrix} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & \gamma_2 & 1 & \\ & \beta_2 & & 1 \end{bmatrix} P_3 = L'_2 P_3,$$

$$P_3 L'_1 = \Pi_{34} L'_1 = \begin{bmatrix} 1 & & & \\ \gamma_1 & 1 & & \\ \alpha_1 & & & 1 \\ \beta_1 & & 1 & \end{bmatrix} = \begin{bmatrix} 1 & & & \\ \gamma_1 & 1 & & \\ \alpha_1 & & 1 & \\ \beta_1 & & & 1 \end{bmatrix} P_3 = L''_1 P_3.$$

Собирая все вместе, приходим к тождеству

$$L_3 P_3 L_2 P_2 L_1 P_1 A = L_3 L'_2 L''_1 P_3 P_2 P_1 A = U,$$

откуда получаем

$$PA = LU,$$

где

$$P = P_3 P_2 P_1, \quad L = (L''_1)^{-1} (L'_2)^{-1} L_3^{-1}.$$

▷₂ Выразите матрицы L''_1 и L'_2 через L_1 , L_2 , P_2 и P_3 .

Благодаря специальной структуре множителей, по аналогии с (2.13), матрицу L легко вычислить:

$$L = \begin{bmatrix} 1 & & & \\ -\gamma_1 & 1 & & \\ -\alpha_1 & -\gamma_2 & 1 & \\ -\beta_1 & -\beta_2 & -\gamma_3 & 1 \end{bmatrix}.$$

В общем случае матрица L может быть представлена в виде (см. (2.14))

$$L = nI - L_{n-1} - L'_{n-2} - L''_{n-3} - \dots - L_1^{(n-2)}. \quad (2.15)$$

Здесь верхний индекс в скобках обозначает количество штрихов:

$$L_k^{(q)} = L_k^{\dots'} = P_{k+q} L_k^{(q-1)} P_{k+q}, \quad 1 \leq k \leq n-2, \quad q \geq 1.$$

Приведенные выше рассуждения можно провести для произвольного n . В результате процесс прямого хода метода Гаусса с выбором главного элемента по столбцу примет вид

$$L_{n-1} L'_{n-2} \dots L_2^{(n-3)} L_1^{(n-2)} P_{n-1} \dots P_2 P_1 A = U.$$

Таким образом по индукции доказывается следующая теорема.

Теорема 2.4 (*LU-разложение с выбором главного элемента*). Пусть $\det A \neq 0$. Тогда матрица A представима в виде

$$PA = LU, \quad (2.16)$$

где P — матрица перестановки, L — нижнетреугольная матрица с единичной главной диагональю, U — верхнетреугольная матрица с ненулевыми элементами на главной диагонали.

▷₃ Докажите данную теорему.

Практическая реализация

Для построения алгоритма вычисления LU -разложения с выбором главного элемента, т. е. разложения (2.16) достаточно внести лишь небольшие изменения в базовый алгоритм на стр. 30. Главное наблюдение, необходимое для этого, уже было сделано выше: матрицы $L_k^{(q)}$, по которым строится матрица L (2.14), представляют собой матрицу L_k , к k -му столбцу которой (и только к нему!) последовательно применены перестановки $P_{k+1}, P_{k+2}, \dots, P_{n-1}$. Учитывая, что элементы матрицы L хранятся на месте обнуляемых элементов исходной матрицы A , указанные преобразования P_q будут применяться к L_k автоматически при перестановке строк матрицы A .

Что же касается итоговой матрицы перестановки

$$P = P_{n-1} \dots P_2 P_1,$$

то нет необходимости хранить ее в явном виде. Достаточно вместо этого строить в процессе разложения вектор перестановок p . Он получается применением перестановок, осуществляемых в ходе алгоритма, к вектору $(1, 2, \dots, n)^T$.

Если разложение (2.16) построено, то решение системы $Ax = b$ осуществляется по схеме

$$PAx = LUx = Pb = b',$$

то есть сначала нужно применить к вектору b перестановку P , переставив его компоненты в соответствии с вектором p :

$$b'_i = b_{p_i}, \quad i = \overline{1, n}.$$

После этого, как обычно, из уравнения $Ly = b'$ прямой подстановкой находится вектор y , а затем из $Ux = y$ обратной подстановкой находится искомое решение x .

▷₄ Постройте полный алгоритм LU -разложения с выбором главного элемента.

2.3. Разложение Холецкого и его друзья

Системы линейных уравнений с симметричными матрицами достаточно часто возникают на практике. В этой лекции мы изучим базовые алгоритмы решения таких систем.

2.3.1. LDL^T -разложение

Теорема 2.5 (LDL^T -разложение). Пусть $A = A^T$ — симметричная вещественная матрица. Если все ее главные угловые миноры отличны от нуля, то существует разложение

$$A = LDL^T, \quad (2.17)$$

где L — нижнетреугольная матрица с единичной главной диагональю, D — диагональная матрица.

Доказательство. По условию существует LU -разложение

$$A = LU = U^T L^T = A^T,$$

откуда $U = L^{-1}U^T L^T$. Обозначим $D = L^{-1}U^T$, тогда $U = DL^T$ и

$$A = LU = LDL^T.$$

Рассмотрим подробнее матрицу $D = L^{-1}U^T$. С одной стороны, D — нижнетреугольная, так как является произведением нижнетреугольных матриц L^{-1} и U^T . С другой стороны, $D = U(L^T)^{-1}$, т. е. D является еще и верхнетреугольной. Следовательно, D — диагональная матрица, что и требовалось доказать. \square

Из теоремы 2.5 получаем, что LDL^T -разложение симметричной матрицы можно найти, просто выполнив ее LU -разложение. Так как $U = DL^T$, а у L единичная главная диагональ, матрицу D легко найти:

$$D = \text{diag}\{u_{11}, u_{22}, \dots, u_{nn}\}.$$

Можно поступить и немного по-другому: просто найти матрицу U прямым ходом метода Гаусса (без построения L). Далее имеем

$$U = DL^T \Rightarrow L = (D^{-1}U)^T, \quad (2.18)$$

то есть матрица L^T находится путем деления каждой строки U на ее диагональный элемент.

Однако, оба этих способа не имеют большого смысла если для выполнения прямого хода метода Гаусса (для построения LU -разложения) мы будем использовать стандартный алгоритм. Самое главное здесь — использовать свойство симметрии матрицы A , так как *прямой ход метода Гаусса можно выполнить в два раза быстрее, если матрица симметрична*. Чтобы в этом убедиться, рассмотрим подробно первый шаг метода для матрицы A . Элементы матрицы $A^{(2)}$, которая получится после этого шага, обозначим a'_{ij} . Тогда по определению для всех i, j от 2 до n будем иметь

$$a'_{ij} = a_{ij} - \frac{a_{i1}}{a_{11}}a_{1j}.$$

Аналогично, учитывая симметрию матрицы A , получим

$$a'_{ji} = a_{ji} - \frac{a_{j1}}{a_{11}}a_{1i} = a_{ij} - \frac{a_{1j}}{a_{11}}a_{i1} = a'_{ij}. \quad (2.19)$$

То есть, после первого шага метода Гаусса матрица системы, не считая первые столбец и строку, остается симметричной! Понятно, что по индукции это утверждение будет выполняться и на всех последующих шагах метода Гаусса. Следовательно, формула (2.19) дает способ практически вдвое сократить количество арифметических операций при построении LU -разложения симметричной матрицы: на каждом шаге метода можно вычислять только элементы верхнего (или нижнего) треугольника, включая диагональ. Один из возможных вариантов алгоритма приведен ниже. После его выполнения на месте верхнего треугольника матрицы A будет находиться матрица L^T (без диагонали), а на диагонали — элементы матрицы D . Обратите внимание на минимальные отличия от метода Гаусса.

```

1: for  $k = \overline{1, n-1}$  do
2:   for  $i = \overline{k+1, n}$  do
3:      $a_{i,i:n} -= (a_{ki}/a_{kk})a_{k,i:n}$ 
4:   end for
5:    $a_{k,k+1:n} /= a_{kk}$ 
6: end for
```

▷₁ Постройте аналогичный алгоритм, который заменяет нижний треугольник A на элементы матрицы L . Какой из этих двух алгоритмов более эффективен при реализации на компьютере? Почему?

▷₂ Запишите схему решения СЛАУ $Ax = b$ если известно LDL^T -разложение матрицы A .

О вычислительной устойчивости метода. Так как при выполнении LDL^T -разложения мы не выполняем выбора главного элемента, рассмотренный алгоритм не может быть применим, если матрица A имеет нулевые главные миноры. Более того, если при $\det[A]_k \approx 0$ метод формально пригоден, то по факту в этом случае обусловленность итоговой матрицы L может быть слишком большой. Поэтому LDL^T -разложение является вычислительно неустойчивым методом решения СЛАУ и применять его в машинной арифметике следует с осторожностью.

2.3.2. Метод квадратного корня

Методом квадратного корня называют метод решения СЛАУ с симметричной и положительно определенной матрицей, основанный на разложении Холецкого. Прежде, чем изучить это разложение, рассмотрим еще раз разложение $A = LDL^T$. Диагональную матрицу D можно представить в виде

$$D = H E H,$$

где H и E — диагональные матрицы, элементы которых равны

$$h_{ii} = \sqrt{|d_{ii}|}, \quad e_{ii} = \text{sign } d_{ii}. \quad (2.20)$$

В силу невырожденности матрицы D такое разложение на множители всегда осуществимо. Тогда

$$A = LDL^T = L H E H^T L^T = G E G^T,$$

где

$$G = L H.$$

Таким образом, следствием теоремы 2.5 является

Теорема 2.6. *Если все главные угловые миноры симметричной вещественной матрицы A отличны от нуля, то существует разложение*

$$A = G E G^T, \quad (2.21)$$

где G — нижнетреугольная матрица, E — диагональная матрица, элементы которой по модулю равны 1.

▷₃ Опишите схему решения системы $Ax = b$, если известно разложение $A = GEG^T$.
Так как $G = LH$, элементы матрицы G могут быть вычислены как

$$g_{ij} = l_{ij} \sqrt{|d_{jj}|}.$$

С другой стороны, если известна матрица $U = DL^T$, то из тождеств $G^T = HL^T$ и (2.18) получим $G^T = HD^{-1}U = E^{-1}H^{-1}U$, то есть

$$g_{ji} = e_{ii}^{-1} h_{ii}^{-1} u_{ij} = \frac{\text{sign } u_{ii}}{\sqrt{|u_{ii}|}} u_{ij}. \quad (2.22)$$

▷₄ Запишите алгоритм построения GEG^T -разложения на основе прямого хода метода Гаусса в двух вариантах: а) заменяя верхний треугольник матрицы A на матрицу G^T , б) заменяя нижний треугольник A на G .

Симметричные положительно определенные матрицы

Вспомним теперь некоторые факты из линейной алгебры.

Определение 2.3. Квадратная вещественная матрица A называется *положительно определенной* ($A > 0$), если

$$(Ax, x) > 0 \quad \forall x \in \mathbb{R}^n, x \neq 0.$$

В дальнейшем нам понадобятся следующие свойства положительно определенных матриц:

- 1) $A^T = A > 0$ тогда и только тогда, когда все ее главные угловые миноры положительны (критерий Сильвестра);
- 2) $A^T = A > 0$ тогда и только тогда, когда все собственные значения A вещественны и положительны.

В частности, по критерию Сильвестра все главные угловые миноры симметричной положительно определенной матрицы отличны от нуля, то есть у таких матриц всегда существуют LDL^T - и GEG^T -разложения (метод Гаусса без выбора главного элемента всегда осуществим). Кроме этого, положительная определенность матрицы A благотворно влияет на вычислительную устойчивость, так что метод квадратного корня (разложение Холецкого), в отличие от LDL^T -разложения, не обладает повышенной чувствительностью к погрешностям округления.

Разложение Холецкого

Теорема 2.7 (разложение Холецкого). Для любой положительно определенной вещественной матрицы A справедливо разложение

$$A = GG^T, \quad (2.23)$$

где G — некоторая нижнетреугольная матрица.

Доказательство. По критерию Сильвестра имеем $\det[A]_k > 0$, $\forall k = \overline{1, n}$, то есть существует разложение $A = GEG^T$. Нетрудно заметить, что $[A]_k = [G]_k[E]_k[G^T]_k$, откуда

$$\det[A]_k = (\det[G]_k)^2 \det[E]_k > 0,$$

следовательно, $\det[E]_k = 1$ для всех $k = \overline{1, n}$, то есть $E = I$. □

Из теоремы 2.7 следует, что в случае $A = A^T > 0$ формула (2.22) превращается просто в

$$g_{ji} = \frac{1}{\sqrt{u_{ii}}} u_{ij},$$

то есть

для вычисления матрицы G^T нужно просто разделить каждую строку матрицы U (которая получается в результате прямого хода метода Гаусса) на корень из диагонального элемента в этой строке.

Положительная определенность матрицы A при этом гарантирует, что $u_{ii} = d_{ii} > 0$. Естественно, прямой ход метода Гаусса следует выполнять рационально, с учетом симметрии матрицы, как мы делали это ранее.

▷₅ Запишите алгоритм построения разложения Холецкого на основе (оптимизированного) прямого хода метода Гаусса в двух вариантах: а) заменяя верхний треугольник матрицы A на матрицу G^T , б) заменяя нижний треугольник A на G .

▷₆ Опишите схему решения системы $Ax = b$, если известно разложение Холецкого $A = GG^T$.

2.4. Метод отражений

2.4.1. Ортогональные преобразования

Определение 2.4. Линейное преобразование $Q : \mathbb{R}^n \rightarrow \mathbb{R}^n$ называется *ортогональным*, если оно сохраняет евклидову норму векторов:

$$\|Qx\|_2 = \|x\|_2 \quad \forall x \in \mathbb{R}^n.$$

Матрица ортогонального преобразования называется ортогональной.

Возможны также следующие эквивалентные определения ортогонального преобразования:

- преобразование ортогонально тогда и только тогда, когда оно сохраняет скалярное произведение:

$$(Qx, Qy) = (x, y) \quad \forall x, y \in \mathbb{R}^n;$$

- преобразование ортогонально тогда и только тогда, когда его матрица удовлетворяет условию

$$Q^{-1} = Q^T.$$

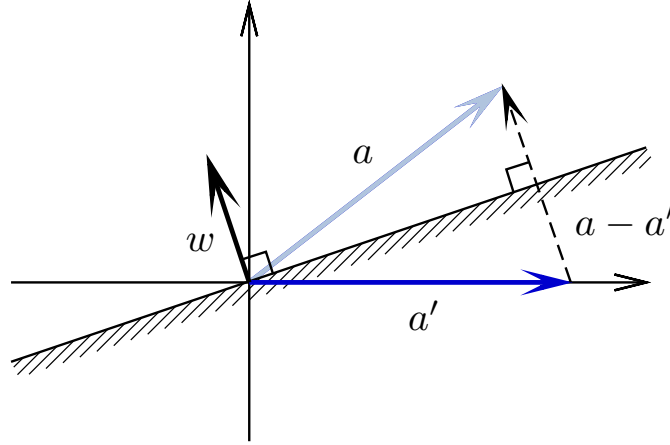
Последняя формулировка наиболее часто приводится в учебниках в качестве определения ортогональных матриц.

Свойства ортогональных матриц

- 1) Строки и столбцы ортогональной матрицы образуют ортонормированные системы векторов.
- 2) Определитель ортогональной матрицы по модулю равен 1.
- 3) Число обусловленности ортогональной матрицы в спектральной норме равно 1.

2.4.2. Преобразование отражения

Важным представителем класса ортогональных преобразований является преобразование отражения, с которым мы сейчас познакомимся. Рассмотрим n -мерном в евклидовом пространстве вектор w единичной длины: $\|w\|_2 = 1$. Этот вектор определяет гиперплоскость W — множество векторов, ортогональных ему. Любой вектор $a \in \mathbb{R}^n$ можно отразить относительно этой гиперплоскости. Используя геометрическую интуицию в двух- или трехмерном случае, можно записать условие, которому должен удовлетворять полученный после отражения вектор a' :



$$(a' - a) \perp W, \quad \text{или} \quad a' - a = \gamma w, \quad \gamma \in \mathbb{R}. \quad (2.24)$$

Отсюда получим $a' = a + \gamma w$. Из еще одного очевидного свойства

$$(a + a') \in W, \quad \text{или} \quad (a + a', w) = 0$$

легко найти неизвестный коэффициент γ получить правило вычисления отраженного вектора:

$$a' = a - 2(a, w)w. \quad (2.25)$$

Легко видеть, что преобразование отражения линейно, то есть задается некоторой матрицей. Найдем вид этой матрицы:

$$a' = a - 2(a, w)w = Ia - 2w(w^T a) = Ia - 2(ww^T)a = (I - 2ww^T)a,$$

то есть (2.25) эквивалентно $a' = Ha$, где H — матрица отражения, однозначно определяемая вектором нормали w по формуле

$$H = I - 2ww^T. \quad (2.26)$$

Преобразование называют также преобразованием Хаусхолдера (Householder).

Замечание 2.2. Заметим, что вычисление преобразования отражения по формуле (2.25) требует $O(n)$ операций умножения и сложения, в то время как умножение на матрицу H требует $O(n^2)$ операций. Поэтому в явном виде матрицу отражения практически никогда не используют, а хранят только w .

▷₁ Докажите, что а) H ортогональна, б) $H^{-1} = H$.

2.4.3. Метод отражений для решения СЛАУ

Рассмотрим систему $Ax = b$. Если для приведения матрицы A к верхнетреугольному виду вместо элементарных преобразований использовать преобразования отражения, получим новый метод решения СЛАУ — метод отражений.

Как и ранее, на k -м шаге прямого хода будут обнуляться все элементы k -го столбца, расположенные ниже главной диагонали. Матрицу системы перед выполнением k -го шага будем обозначать $A^{(k)}$. В итоге мы должны получить

$$Q_{n-1} \dots Q_2 Q_1 A = R,$$

где Q_k — ортогональные матрицы, определяющие преобразования отражения, а R — верхнетреугольная матрица, аналог матрицы U в LU -разложении.

Рассмотрим подробно первый шаг, для краткости обозначая $\|\cdot\| = \|\cdot\|_2$. Пусть a_1 — первый столбец матрицы A . После первого шага этот вектор должен быть преобразован в

$$a'_1 = (a'_{11}, 0, \dots, 0)^T.$$

Так как это преобразование является ортогональным, имеем $\|a_1\| = \|a'_1\|$, откуда

$$a'_{11} = \pm \|a_1\|. \quad (2.27)$$

Теперь необходимо найти преобразование отражения, которое отображает a_1 в a'_1 . Соответствующий вектор нормали w_1 можно легко найти из (2.24). Учитывая, что норма этого вектора должна быть единичной, будем иметь

$$w_1 = \|a_1 - a'_1\|^{-1}(a_1 - a'_1). \quad (2.28)$$

Для того, чтобы избежать вычитания близких чисел при использовании этой формулы, следует воспользоваться свободой выбора знака для a'_{11} в формуле (2.27). Чтобы операция вычитания на самом деле представляла собой сложение одинаковых по модулю чисел, следует выбрать

$$a'_{11} = (-\operatorname{sign} a_{11})\|a_1\|.$$

Итак, мы нашли преобразование отражения, которое приводит первый столбец матрицы A к нужному виду: $Q_1 a_1 = a'_1$, где Q_1 — матрица отражения $Q_1 = I - 2w_1 w_1^T$, которую нет смысла вычислять в явном виде.

Чтобы получить $A^{(2)} = Q_1 A$ нам нужно применить Q_1 к остальным столбцам a_j матрицы A , что можно сделать по формуле (2.25):

$$a'_j = a_j - 2(a_j, w_1)w_1, \quad j = \overline{2, n}.$$

Не забудем также и про вектор правой части:

$$b' = b^{(2)} = b - 2(b, w_1)w_1.$$

Первый шаг метода отражений выполнен.

Следующие шаги выполняются аналогично: на k -ом шаге описанные выше преобразования применяются к подматрице, составленной из последних $n - k + 1$ строк и столбцов матрицы $A^{(k)}$. Соответствующее преобразование $A^{(k+1)} = Q_k A^{(k)}$ задается блочной матрицей

$$Q_k = \begin{bmatrix} I_{k-1} & 0 \\ 0 & H_k \end{bmatrix},$$

где H_k — матрица отражения, определяемая вектором нормали $w_k \in \mathbb{R}^{n-k+1}$, который вычисляется аналогично (2.28). Как видно, на каждом шаге в столбцах матрицы $A^{(k)}$ и в векторе b под действием преобразований Q_k будут меняться компоненты с индексами от k до n . Еще раз подчеркнем, что матрицы H_k и Q_k в явном виде не строятся, а все преобразования отражения выполняются по правилу (2.25).

▷₂ Запишите алгоритм прямого хода метода отражений для решения СЛАУ.

По окончании прямого хода останется обратной подстановкой решить систему $Rx = b^{(n)}$, где $R = A^{(n)}$ — итоговая верхнетреугольная матрица, а $b^{(n)}$ — преобразованный вектор правой части:

$$b^{(n)} = Q_{n-1} \dots Q_1 b.$$

Отметим, что $b^{(n)}$ можно вычислить и после прямого хода, если сохранены все векторы нормали w_k , которые строились во время прямого хода. В этом случае можно назвать описанный подход решением СЛАУ методом QR -разложения (см. далее).

Трудоемкость прямого хода метода отражений в два раза больше, чем у метода Гаусса.

▷₃ Докажите это утверждение.

▷₄ Как в ходе метода отражений распознать случай вырожденной матрицы A ?

2.4.4. QR-разложение

Определение 2.5. Пусть A — матрица размера $m \times n$, причем $m \geq n$. *QR-разложением* такой матрицы называется представление ее в виде

$$A = QR,$$

где Q — ортогональная матрица ($Q^{-1} = Q^T$) размера $m \times m$, R — верхнетреугольная матрица ($r_{ij} = 0 \quad \forall i > j$), размер которой совпадает с размером матрицы A .

Приведем пример для случая $m = 5, n = 3$:

$$\underbrace{\begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}}_A = \underbrace{\begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}}_Q \underbrace{\begin{bmatrix} \times & \times & \times \\ & \times & \times \\ & & \times \\ & & & \\ & & & & \end{bmatrix}}_R$$

Заметим, что при $m > n$ последние $m - n$ столбцов матрицы Q никак не влияют на элементы матрицы A (так как последние $m - n$ строк матрицы R равны нулю), поэтому часто в таких случаях эти лишние столбцы и строки отбрасывают и рассматривают так называемое *сокращенное QR-разложение* $A = \hat{Q}\hat{R}$:

$$\underbrace{\begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}}_A = \underbrace{\begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}}_{\hat{Q}} \underbrace{\begin{bmatrix} \times & \times & \times \\ & \times & \times \\ & & \times \\ & & & \\ & & & & \end{bmatrix}}_{\hat{R}}.$$

В общем случае размеры матриц \hat{Q} и \hat{R} равны $m \times n$ и $n \times n$ соответственно. Матрицу \hat{Q} уже нельзя считать ортогональной, но у нее остается важное свойство ортогональных матриц: все ее n столбцов образуют ортонормированную систему, то есть $\hat{Q}^T \hat{Q} = I$.

Построение QR-разложения методом отражений

Легко увидеть, что в случае квадратной матрицы прямой ход метода отражений, описанный в пункте 2.4.3, представляет собой процесс

построения QR -разложения матрицы A :

$$Q_{n-1} \dots Q_2 Q_1 A = R,$$

откуда получаем $A = QR$, где R — верхнетреугольная матрица, а ортогональная матрица Q имеет вид

$$Q = (Q_{n-1} \dots Q_1)^{-1} = Q_1 \dots Q_{n-1},$$

так как $Q_k^{-1} = Q_k^T = Q_k$.

При $m > n$ суть метода остается той же, только для приведения матрицы A к верхнетреугольному виду нужно будет n (а не $n - 1$) преобразований отражения. Рассмотрим схему построения разложения для $m = 5, n = 3$:

$$\underbrace{\begin{bmatrix} \otimes & \times & \times \\ \otimes & \times & \times \\ \otimes & \times & \times \\ \otimes & \times & \times \\ \otimes & \times & \times \end{bmatrix}}_{A^{(1)}=A} \xrightarrow{Q_1} \underbrace{\begin{bmatrix} \times & \times & \times \\ & \otimes & \times \\ & \otimes & \times \\ & \otimes & \times \\ & \otimes & \times \end{bmatrix}}_{A^{(2)}=Q_1 A^{(1)}} \xrightarrow{Q_2} \underbrace{\begin{bmatrix} \times & \times & \times \\ & \times & \times \\ & & \otimes \\ & & \otimes \\ & & \otimes \end{bmatrix}}_{A^{(3)}=Q_2 A^{(2)}} \xrightarrow{Q_3} \underbrace{\begin{bmatrix} \times & \times & \times \\ & \times & \times \\ & & \times \end{bmatrix}}_{R=Q_3 Q_2 Q_1 A}$$

Кружками на этой схеме выделены компоненты вектора, который на соответствующем этапе играет роль вектора a_1 , по которому строится нормаль (2.28), определяющая преобразование отражения.

В общем случае получится

$$\underbrace{Q_n \dots Q_2 Q_1}_{Q^{-1}} A = R,$$

откуда

$$Q = (Q_n \dots Q_2 Q_1)^{-1} = Q_1 Q_2 \dots Q_n, \quad (2.29)$$

$$Q_k = \begin{bmatrix} I_{k-1} & 0 \\ 0 & H_k \end{bmatrix},$$

где I_{k-1} — единичная матрица размера $k - 1$, H_k — матрица отражения размера $m - k + 1$, R — верхнетреугольная матрица.

Теперь вспомним, что вместо матрицы Q проще хранить векторы нормали w_1, \dots, w_n , которые однозначным образом определяют матрицы $H_k = I - 2w_k w_k^T$ и, следовательно, «сборную» матрицу Q по формуле

(2.29). В силу блочной структуры матрицы Q_k при умножении на вектор она изменяет в нем лишь координаты с индексами от k до m .

На практике обычно требуется вычислять преобразование Q^T ,

$$Q^T = Q^{-1} = Q_n \dots Q_2 Q_1.$$

Если использовать стандартное обозначение $x_{k:m} = (x_k, x_{k+1}, \dots, x_m)^T$, алгоритм вычисления $x \leftarrow Q^T x$ можно записать в следующем виде:

- 1: **for** $k = \overline{1, n}$ **do**
- 2: $x_{k:m} \leftarrow x_{k:m} - 2(x_{k:m}, w_k)w_k$
- 3: **end for**

2.5. Задача наименьших квадратов

2.5.1. Постановка задачи

Рассмотрим переопределенную систему линейных уравнений, то есть систему, в которой количество неизвестных n меньше количества уравнений m :

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}. \quad (2.30)$$

Матрица системы A является прямоугольной, строк в ней больше чем столбцов. В дальнейшем будем предполагать, что ранг этой матрицы равен n , то есть все ее столбцы линейно независимы. Решение системы (2.30) существует лишь в случае, если вектор $b \in \mathbb{R}^m$ принадлежит линейной оболочке системы векторов-столбцов матрицы A :

$$b \in V = \text{span}\{a_1, \dots, a_n\}.$$

Чтобы решение существовало при всех b нужно придать уравнениям (2.30) несколько иной смысл, а именно: будем искать вектор $x \in \mathbb{R}^n$ такой, который минимизировал бы норму невязки $Ax - b$:

$$\|Ax - b\|_2 \rightarrow \min. \quad (2.31)$$

Задачи такого рода часто возникают в приложениях. Например, пусть задан набор точек на плоскости (ξ_i, η_i) , $i = \overline{1, m}$. Требуется найти многочлен φ степени $n - 1$, график которого был бы «как можно более близок» к этим точкам. Это условие можно сформулировать как

$$\sum_{i=1}^m (\varphi(\xi_i) - \eta_i)^2 \rightarrow \min.$$

Подставляя сюда $\varphi(\xi) = \sum_{j=1}^n c_j \xi^{j-1}$, получаем линейную задачу наименьших квадратов для нахождения коэффициентов (c_j) :

$$\|Ac - b\|_2 \rightarrow \min,$$

где $a_{ij} = \xi_i^{j-1}$, $i = \overline{1, m}$, $j = \overline{1, n}$, $c = (c_1, \dots, c_n)^T$, $b = (\eta_1, \dots, \eta_m)^T$.

2.5.2. Решение с помощью нормальных уравнений

Определение 2.6. Ортогональной проекцией вектора $x \in \mathbb{R}^m$ на подпространство $V \subset \mathbb{R}^m$ называется вектор $y \in V$ такой, что $x - y \perp v \quad \forall v \in V$.

Если известен базис $\{a_1, \dots, a_n\}$ подпространства V , то эквивалентным определением можно считать следующее: вектор $y \in V$ является ортогональной проекцией вектора x , если

$$x - y \perp a_i \quad \forall i = \overline{1, n}. \quad (2.32)$$

▷₁ Докажите эквивалентность этих двух определений.

Основное свойство ортогональной проекции вектора x состоит в том, что она является наилучшим среднеквадратичным приближением к x среди всех элементов подпространства V :

$$\|x - y\|_2 = \min_{v \in V} \|x - v\|_2.$$

Этот факт, вообще говоря, следует доказывать, но мы этого делать не будем, полагаясь на геометрическую интуицию.

Исходя из сказанного, решение задачи наименьших квадратов (2.31) можно найти из условий (2.32). Действительно, по условию вектор Ax должен быть проекцией вектора b на $V = \text{span}\{a_1, \dots, a_n\}$, то есть невязка $Ax - b$ должна быть ортогональна каждому базисному вектору подпространства V . Это условие записывается в виде

$$(Ax - b, a_j) = 0, \quad \text{или} \quad (Ax, a_j) = (b, a_j) \quad \forall j = \overline{1, n}.$$

Полученные условия называются *нормальными уравнениями*. Записывая их для всех j нетрудно заметить, что они эквивалентны СЛАУ

$$A^T Ax = A^T b. \quad (2.33)$$

Таким образом, решение задачи наименьших квадратов (2.31) совпадает с решением СЛАУ (2.33). Матрица $A^T A$ является симметричной и положительно определенной, поэтому для решения нормальных уравнений лучше использовать метод квадратного корня (разложение Холецкого). Если исходная матрица A имеет ранг n , то матрица $A^T A$ невырождена, и, следовательно, решение задачи наименьших квадратов существует и единственно.

▷₂ Подсчитайте количество операций, необходимое для нахождения x из уравнения (2.33).

Решение задачи наименьших квадратов путем сведения к нормальным уравнениям является самым экономичным, но, в то же время, и наименее вычислительно устойчивым способом. Более надежный метод основан на QR -разложении матрицы A .

2.5.3. Решение с помощью QR -разложения

Рассмотрим «нормальную» систему $A^T A x = A^T b$ и сокращенное QR -разложение $A = \hat{Q} \hat{R}$:

$$A^T A x = A^T b \quad \Leftrightarrow \quad \hat{R}^T \hat{Q}^T \hat{Q} \hat{R} x = \hat{R}^T \hat{Q}^T b,$$

откуда

$$\hat{R} x = \hat{Q}^T b. \quad (2.34)$$

Отсюда получаем следующий алгоритм решения задачи наименьших квадратов:

- 1) Построить сокращенное QR -разложение $A = \hat{Q} \hat{R}$.
- 2) Вычислить вектор $\hat{Q}^T b$.
- 3) Решить СЛАУ (2.34) обратной подстановкой.

Как было сказано в п. 2.4.4, матрица Q как правило не хранится в явном виде, а определяется набором векторов нормали $\{w_k\}_{k=1}^n$ по формулам (2.29). В связи с этим обсудим схему вычисления вектора $b' = \hat{Q}^T b$ в пункте 2. Матрица \hat{Q}^T это первые n строк матрицы Q^T , поэтому вектор b' содержит первые n компонент вектора $Q^T b$. Значит, сначала нужно найти

$$b^{(n)} = Q_n \dots Q_2 Q_1 b,$$

см. алгоритм в конце п. 2.4.4, после чего положить

$$b' = b_{1:n}^{(n)}.$$

2.6. Форматы хранения разреженных матриц

Определение 2.7. *Разреженными* называют матрицы, содержащие большой процент нулевых элементов.

Разреженные матрицы очень часто возникают в приложениях, в частности, при численном моделировании различных физических процессов. Количество ненулевых элементов матрицы в дальнейшем будем обозначать n_z .

Хранить в памяти ЭВМ все нулевые элементы разреженных матриц нерационально. Они, во-первых, зря занимают память, и, во-вторых, замедляют операции с матрицами. Поэтому существует ряд общепринятых способов хранения разреженных матриц.

2.6.1. Координатный формат

Самый простой способ — так называемый *координатный формат*. В этом формате для хранения вещественной матрицы A используется три массива:

- AA — массив вещественных чисел для хранения ненулевых элементов матрицы A ;
- IA — массив целых чисел для хранения номеров строк соответствующих элементов массива AA ;
- JA — аналогичный массив для номеров столбцов.

Длина всех трех массивов равна n_z .

Пример 2.1. Записать в координатном формате матрицу

2.				
		5.	7.	
		9.		
1.	4.		3.	
	8.			6.

Решение. Представление матрицы в координатном формате задается с точностью до перестановки элементов:

AA	1.	2.	3.	4.	5.	6.	7.	8.	9.
IA	4	1	4	4	2	5	2	5	3
JA	1	1	4	2	3	5	4	2	3



2.6.2. Форматы CSR и CSC

Если в примере 2.1 упорядочить элементы матрицы построчно, то информация в массиве IA окажется избыточной:

AA	2.	5.	7.	9.	1.	4.	3.	8.	6.
IA	1	2	2	3	4	4	4	5	5
JA	1	3	4	3	1	2	4	2	5

В этом случае достаточно хранить лишь *указатель* на элемент массива AA, с которого начинается хранение i -й строки:

IA	1	2	4	5	8	10
----	---	---	---	---	---	----

Здесь последний элемент вектора IA служит для того, чтобы можно было определить, где заканчивается n -я строка. Описанный способ хранения матриц называется *форматом CSR* (Compressed Sparse Row).

Если ненулевые элементы в массиве AA упорядочить не по строкам, а по столбцам, по аналогии получится *формат CSC* (Compressed Sparse Column).

Описание формата CSR (CSC)

- AA — массив вещественных чисел длины n_z , в котором хранятся упорядоченные по строкам (по столбцам) ненулевые элементы A .
- JA — массив целых чисел длины n_z для хранения номеров столбцов (строк) соответствующих элементов массива AA.
- IA — массив целых чисел длины $n + 1$. На i -й позиции массива хранится *номер позиции* в массиве AA, с которой начинается хранение элементов i -й строки (столбца). Более конкретно: $IA[1]=1$, $IA[i+1]=IA[i]+n_i$, где n_i — число ненулевых элементов в i -й строке (столбце).

Формат CSR является одним из наиболее популярных.

▷₁ Укажите достоинства и недостатки этого формата по сравнению с координатным форматом.

Алгоритм умножения матрицы в формате CSR на вектор

Вход: n , AA, JA, IA, x .

Выход: y .

```

for i=1 to n do
  y[i]=0
  for j=IA[i] to IA[i+1]-1 do
    y[i]=y[i]+AA[j]*x[JA[j]]
  end for
end for

```

Алгоритм умножения матрицы в формате CSC на вектор

```

y=0
for j=1 to n do
  for i=IA[j] to IA[j+1]-1 do
    y[?]=y[?]+AA[?]*x[?]
  end for
end for

```

▷₂ Что должно стоять вместо знаков вопроса?

2.6.3. Формат MSR

Во многих итерационных методах решения СЛАУ диагональные элементы матрицы A играют особую роль: доступ к ним нужно осуществлять чаще, чем к другим элементам матрицы. Рассмотренные ранее форматы не позволяют быстро найти диагональные элементы. Решить эту проблему помогает модификация формата CSR — *формат MSR* (Modified Sparse Row).

Основные отличия этого формата от формата CSR состоят в следующем:

- 1) изменен формат массива AA: сначала в него *полностью* записывается главная диагональ матрицы A , а затем — недиагональные ненулевые элементы;
- 2) массивы IA и JA объединены в один (назовем его IJ).

Рассмотрим матрицу из примера 2.1:

2.				
		5.	7.	
		9.		
1.	4.		3.	
	8.			6.

Запишем ее представление в формате MSR:

	1	2	3	4	5	6	7	8	9	10	11
AA	2.	0.	9.	3.	6.	×	5.	7.	1.	4.	8.
IJ	7	7	9	9	11	12	3	4	1	2	2

Массив AA: первые n элементов занимает главная диагональ A . Элемент на позиции $n + 1$ *не используется*. Начиная с $(n + 2)$ -го элемента построчно хранятся недиагональные ненулевые элементы A .

Массив IJ: $IJ[1]=n+2$, $IJ[i+1]=IJ[i]+n_i$ для $i = \overline{1, n}$, где n_i — количество недиагональных ненулевых элементов в i -й строке. Далее — номера столбцов для соответствующих элементов массива AA.

Замечание 2.3. Массивы AA и IJ полностью описывают матрицу A . Размерность матрицы легко найти по значению $IJ[1]$:

$$n = IJ[1] - 2.$$

Кроме этого, длина обоих массивов равна $IJ[n+1] - 1$.

Алгоритм умножения матрицы в формате MSR на вектор

Вход: AA, IJ, x.

Выход: y.

```

n=IJ[1]-2
for i=1 to n do
    y[i]=AA[i]*x[i]
    for j=IJ[i] to IJ[i+1]-1 do
        y[i]=y[i]+AA[j]*x[IJ[j]]
    end for
end for

```

▷₃ Почему в массиве AA пустует элемент номер $n + 1$?

2.7. Итерационные методы решения СЛАУ

2.7.1. Общая характеристика

Все методы решения СЛАУ, которые мы до сих пор рассматривали, являлись *прямыми*, или *точными*, т. е. позволяли (при отсутствии ошибок округления) найти точное решение системы $Ax = b$ за конечное число операций. Сейчас мы рассмотрим другой класс методов — *итерационные* методы, которые позволяют за конечное число операций найти решение лишь *приближенно*, но с произвольной наперед заданной точностью¹. Это значит, что в общем случае для нахождения точного решения итерационному методу потребовалось бы *бесконечно много* операций. Однако это не проблема, так как даже точные методы в условиях машинной арифметики всегда дают решение с точностью до ошибок округления.

Если говорить абстрактно, то итерационный метод (процесс) решения СЛАУ может рассматриваться как метод построения последовательности векторов

$$x^0, x^1, \dots, x^k, \dots,$$

такой, что если существует ее предел x^* ,

$$\|x^k - x^*\| \xrightarrow{k \rightarrow \infty} 0,$$

то $x^* = A^{-1}b$ — искомое решение.

2.7.2. Стационарные итерационные процессы

Принцип построения

Рассмотрим СЛАУ

$$Ax = b$$

с невырожденной матрицей A , а также эквивалентную ей систему

$$x = Bx + g, \tag{2.35}$$

¹Естественно, чем выше требуемая точность, тем больше вычислительной работы придется проделать.

где B — матрица; g — вектор соответствующей размерности. Эквивалентность здесь означает, что эти системы имеют одинаковое решение

$$A^{-1}b = (I - B)^{-1}g.$$

Из (2.35) автоматически получается итерационный процесс

$$x^{k+1} = Bx^k + g, \quad k = 0, 1, 2, \dots \quad (2.36)$$

Нетрудно заметить, что если такой процесс сходится, т. е. если сходится последовательность (x^k) , то предел этой последовательности x^* удовлетворяет условию $x^* = Bx^* + g$, откуда следует, что $x^* = A^{-1}b$ — искомое решение. Процессы типа (2.36) будем называть *стационарными итерационными процессами*. Слово «стационарный» здесь обозначает тот факт, что матрица B не зависит от k . Строго говоря, к этому определению стоило бы добавить «одношаговый», так как следующее приближение x^{k+1} зависит только от одного предыдущего x^k , но для краткости мы этого делать не будем.

Рассмотрим наиболее известные примеры таких процессов.

Классические стационарные итерационные методы

Метод Якоби. Рассмотрим i -е уравнение СЛАУ $Ax = b$:

$$a_{i1}x_1 + \dots + a_{ii}x_i + \dots + a_{in}x_n = b_i.$$

Выражая из него x_i , получим

$$x_i = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij}x_j \right), \quad i = \overline{1, n}. \quad (2.37)$$

Для того чтобы записать это тождество в векторном виде, рассмотрим разбиение матрицы A на слагаемые согласно рис. 2.4:

$$A = L + D + R. \quad (2.38)$$

Тогда (2.37) примет вид

$$x = D^{-1}(b - (L + R)x) = B_Jx + g_J,$$

где

$$B_J = -D^{-1}(L + R), \quad g_J = D^{-1}b. \quad (2.39)$$

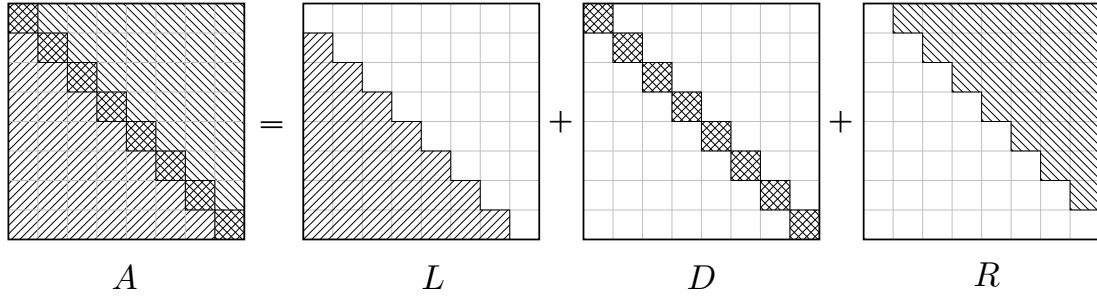


Рис. 2.4

Соответствующий системе (2.37) итерационный метод

$$x_i^{k+1} = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij} x_j^k \right), \quad i = \overline{1, n}, \quad k = 0, 1, 2, \dots \quad (2.40a)$$

называется *методом Якоби*. Его векторная форма имеет вид

$$x^{k+1} = B_J x^k + g_J, \quad (2.40б)$$

где B_J и g_J определяются по формуле (2.39).

Заметим, что $L + R = A - D$, поэтому для матрицы B_J существует альтернативная форма записи

$$B_J = I - D^{-1}A.$$

Метод Гаусса – Зейделя. Рассмотрим i -й шаг k -й итерации метода Якоби:

$$x_i^{k+1} = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij} x_j^k \right).$$

К этому моменту нам уже известны компоненты вектора x^{k+1} с номерами от 1 до $i - 1$. Эти компоненты *могут быть* более точны, чем соответствующие компоненты текущего приближения x^k , поэтому их можно использовать в сумме (2.40a):

$$x_i^{k+1} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{k+1} - \sum_{j=i+1}^n a_{ij} x_j^k \right). \quad (2.41a)$$

Векторный вариант (2.41a) имеет вид

$$x^{k+1} = D^{-1}(b - Lx^{k+1} - Rx^k),$$

откуда

$$\begin{aligned} x^{k+1} &= B_S x^k + g_S, \\ B_S &= -(D + L)^{-1} R, \quad g_S = (D + L)^{-1} b. \end{aligned} \quad (2.41б)$$

Формулы (2.41а), (2.41б) определяют *метод Гаусса – Зейделя*.

▷₁ Постройте модификацию метода Гаусса – Зейделя для случая, когда компоненты вектора x^{k+1} обновляются в обратном порядке (такой метод называется *обратным методом Гаусса – Зейделя*).

Метод релаксации. *Метод релаксации* получается путем взвешенного осреднения текущего приближения и приближения, построенного по методу Гаусса – Зейделя:

$$x_i^{k+1} = (1 - \omega)x_i^k + \frac{\omega}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{k+1} - \sum_{j=i+1}^n a_{ij}x_j^k \right), \quad (2.42а)$$

где ω — весовой коэффициент, обычно $\omega \in (0, 2)$. Формула (2.42а) в векторной форме имеет вид

$$x^{k+1} = (1 - \omega)x^k + \omega D^{-1}(b - Lx^{k+1} - Rx^k).$$

Умножая обе части на D и группируя слагаемые, получаем

$$\begin{aligned} x^{k+1} &= B_\omega x^k + g_\omega, \\ B_\omega &= (D + \omega L)^{-1}((1 - \omega)D - \omega R), \quad g_\omega = (D + \omega L)^{-1}b. \end{aligned} \quad (2.42б)$$

Замечание 2.4. При $\omega = 1$ метод релаксации очевидно превращается в метод Гаусса – Зейделя.

Внимание! Для программной реализации классических итерационных методов используются исключительно их скалярные формы (2.40а), (2.41а), (2.42а). Соответствующие векторные формы записи (2.40б), (2.41б), (2.42б) используются для анализа сходимости методов (см. далее).

2.7.3. Сходимость стационарных итерационных процессов

До сих пор мы строили итерационные процессы формально, ничего не говоря об их сходимости. Для их обоснования необходимо изучить

сходимость стационарных итерационных процессов

$$x^{k+1} = Bx^k + g.$$

Для начала вспомним некоторые определения из линейной алгебры.

Определение 2.8. Пусть A — квадратная матрица. Вектор $x \neq 0$, $x \in \mathbb{C}^n$, называется *собственным вектором* матрицы A , если существует $\lambda \in \mathbb{C}$ такое, что

$$Ax = \lambda x.$$

Число λ называется *собственным значением*, соответствующим x . Множество всех собственных значений A называется *спектром* и обозначается $\sigma(A)$.

Определение 2.9. *Спектральным радиусом* $\rho(A)$ называется величина наибольшего по модулю собственного значения матрицы A :

$$\rho(A) = \max_{\lambda \in \sigma(A)} |\lambda|.$$

Основной результат о сходимости выглядит следующим образом.

Теорема 2.8 (критерий сходимости стационарных итерационных процессов). *Итерационный процесс $x^{k+1} = Bx^k + g$ сходится при любом начальном приближении x^0 тогда и только тогда, когда все собственные значения матрицы B по модулю меньше единицы, т. е. если $\rho(B) < 1$.*

Доказательство этой теоремы основано на следующих двух леммах.

Лемма 2.1. *Стационарный итерационный процесс $x^{k+1} = Bx^k + g$ сходится при любом начальном приближении тогда и только тогда, когда последовательность матриц $I, B, B^2, \dots, B^k, \dots$ сходится к нулевой матрице, т. е. $\|B^k\| \rightarrow 0$ в некоторой подчиненной матричной норме.*

Доказательство. Рассмотрим погрешность на k -й итерации:

$$x^k - x^* = Bx^{k-1} + g - Bx^* - g = B(x^{k-1} - x^*) = \dots = B^k(x^0 - x^*).$$

Отсюда видно, что если процесс сходится, то есть $\|x^* - x^k\| \rightarrow 0$ при любом $x^0 \in \mathbb{R}^n$, тогда $\|B^k x\| \rightarrow 0$ для всех $x \in \mathbb{R}^n$, т. е. $\|B^k\| \rightarrow 0$.

Пусть теперь $B^k \rightarrow 0$. Тогда

$$\|x^k - x^*\| = \|B^k(x^0 - x^*)\| \leq \|B^k\| \|x^0 - x^*\| \rightarrow 0,$$

при любом $x^0 \in \mathbb{R}^n$. □

Определение 2.10. Матрица A называется *диагонализируемой*, или *матрицей простой структуры*, если она представима в виде

$$A = X\Lambda X^{-1},$$

где матрица X невырождена, а

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n).$$

Здесь λ_i являются, очевидно, собственными значениями A , так как преобразование подобия $X\Lambda X^{-1}$ сохраняет спектр.

Лемма 2.2. *Матричная последовательность*

$$I, B, B^2, \dots, B^k, \dots$$

сходится к нулевой матрице тогда и только тогда, когда $\rho(B) < 1$.

Доказательство. Мы докажем эту лемму лишь для частного случая, когда матрица B диагонализируема:

$$B = X\Lambda X^{-1}.$$

В этом случае имеем

$$B^2 = X\Lambda X^{-1}X\Lambda X^{-1} = X\Lambda^2 X^{-1}$$

и вообще

$$B^k = X\Lambda^k X^{-1},$$

откуда

$$B^k \rightarrow 0 \Leftrightarrow \Lambda^k \rightarrow 0 \Leftrightarrow |\lambda_i|^k \rightarrow 0 \Leftrightarrow \rho(B) < 1. \quad \square$$

Собирая вместе две последние леммы, получаем доказательство теоремы 2.8, из которой теперь можно получить и достаточное условие сходимости.

Следствие 2.2 (достаточное условие сходимости). *Если для некоторой подчиненной матричной нормы $\|B\| < 1$, то итерационный процесс $x^{k+1} = Bx^k + g$ сходится.*

Доказательство. Рассмотрим любое собственное значение λ матрицы B и соответствующий ему собственный вектор ξ с единичной нормой. Тогда в любой подчиненной матричной норме имеем

$$\|B\| \geq \|B\xi\| = \|\lambda\xi\| = |\lambda| \|\xi\| = |\lambda|,$$

откуда

$$\rho(B) \leq \|B\| < 1.$$

Значит, итерационный процесс сходится по теореме 2.8. \square

Замечание 2.5. Из доказательства леммы 2.2 можно увидеть, что в случае диагоналируемой матрицы норма погрешности $\|x^k - x^*\|$ стремится к нулю как геометрическая прогрессия со знаменателем $\rho(B)$:

$$\|x^k - x^*\| = \|B^k(x^0 - x^*)\| = \|X\Lambda^k X^{-1}(x^0 - x^*)\| \leq C\rho(B)^k,$$

$C = \kappa(X)\|x^0 - x^*\|$. Это свойство будет справедливым и в общем случае, однако лишь начиная с некоторого достаточно большого номера итерации k .

Сходимость классических итерационных методов

Применим полученные выше общие результаты о сходимости к методу Якоби (2.40).

Лемма 2.3. *Рассмотрим СЛАУ $Ax = b$. Если матрица A имеет строгое диагональное преобладание, то метод Якоби для такой системы сходится при любом начальном приближении.*

Доказательство. Рассмотрим матрицу $B_J = I - D^{-1}A$:

$$B_J = - \begin{bmatrix} 0 & \frac{a_{12}}{a_{11}} & \frac{a_{13}}{a_{11}} & \dots & \frac{a_{1n}}{a_{11}} \\ \frac{a_{21}}{a_{22}} & 0 & \frac{a_{23}}{a_{22}} & \dots & \frac{a_{2n}}{a_{22}} \\ \frac{a_{31}}{a_{33}} & \frac{a_{32}}{a_{33}} & 0 & \dots & \frac{a_{3n}}{a_{33}} \\ \dots & \dots & \dots & \dots & \dots \\ \frac{a_{n1}}{a_{nn}} & \frac{a_{n2}}{a_{nn}} & \frac{a_{n3}}{a_{nn}} & \dots & 0 \end{bmatrix}.$$

Вычислим максимум-норму этой матрицы:

$$\|B_J\|_\infty = \max_i \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|} = \max_i \frac{\sum_{j \neq i} |a_{ij}|}{|a_{ii}|}.$$

По условию матрица A имеет строгое диагональное преобладание, т. е.

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \forall i = \overline{1, n}.$$

Отсюда сразу следует, что

$$\|B_J\|_\infty < 1,$$

т. е. метод Якоби сходится по следствию 2.2. □

Замечание 2.6. Аналогичное утверждение можно доказать и для метода Гаусса — Зейделя.

Следующий важный результат мы приведем без доказательства.

Теорема 2.9. *Если матрица A является симметричной и положительно определенной, то метод релаксации (2.42) сходится при любом $\omega \in (0, 2)$.*

Следствие 2.3. *Если матрица A является симметричной и положительно определенной, то метод Гаусса — Зейделя (2.41) сходится.*

2.8. Чебышевское ускорение итерационных процессов

Вспомним некоторые предварительные факты из линейной алгебры. Если p — некоторый многочлен степени m общего вида $p(t) = c_m t^m + c_{m-1} t^{m-1} + \dots + c_1 t + c_0$, то его значения определены не только для вещественных чисел, но и для квадратных матриц:

$$p(A) = c_m A^m + c_{m-1} A^{m-1} + \dots + c_1 A + c_0 I.$$

При этом нетрудно доказать следующее важное свойство:

$$\sigma(p(A)) = \{p(\lambda_1), \dots, p(\lambda_n)\}, \quad \lambda_i \in \sigma(A),$$

то есть собственные значения матрицы $p(A)$ равны значениям многочлена p на спектре матрицы A .

▷₁ Докажите это утверждение.

2.8.1. Общая идея метода

Рассмотрим начальное приближение x^0 и выполним m итераций некоторого стационарного метода:

$$x^1 = Bx^0 + g, \dots, x^m = Bx^{m-1} + g.$$

Будем предполагать что этот процесс сходится и спектр B веществен, то есть

$$\sigma(B) \subset [-\rho, \rho] \subset (-1, 1).$$

Составим линейную комбинацию из полученных приближений:

$$y^m = \sum_{k=0}^m c_{m,k} x^k,$$

и попробуем выбрать коэффициенты $c_{m,k}$ так, чтобы новое приближение y^m было как можно более точным. Для этого рассмотрим погрешность

$$y^m - x^* = \sum_{k=0}^m c_{m,k} x^k - x^* = \sum_{k=0}^m c_{m,k} (x^k - x^*) - \left(1 - \sum_{k=0}^m c_{m,k}\right) x^*. \quad (2.43)$$

Первым условием, которое мы наложим на коэффициенты $c_{m,k}$, будет

$$\sum_{k=0}^m c_{m,k} = 1. \quad (2.44)$$

Во-первых, по определению y^m это условие гарантирует, что если все x^k равны x^* , то $y^m = x^*$. Во-вторых, из (2.43) тогда будем иметь

$$\begin{aligned} \|y^m - x^*\| &= \left\| \sum_{k=0}^m c_{m,k} (x^k - x^*) \right\| = \left\| \sum_{k=0}^m c_{m,k} B^k (x^0 - x^*) \right\| = \\ &= \|p_m(B)(x^0 - x^*)\| \leq \|p_m(B)\| \|x^0 - x^*\|, \end{aligned} \quad (2.45)$$

где

$$p_m(t) = \sum_{k=0}^m c_{m,k} t^k.$$

Как и ранее, предположим, что матрица B диагонализируема, $B = X \Lambda X^{-1}$, тогда

$$\|p_m(B)\| \leq \|X\| \|X^{-1}\| \|p_m(\Lambda)\| = \kappa(X) \max_{\lambda \in \sigma(B)} |p_m(\lambda)|.$$

Эта важная оценка наряду с (2.45) говорит о том, что если мы заинтересованы в наименьшей норме погрешности $\|y^m - x^*\|$, то коэффициенты $c_{m,k}$ должны быть такими, *чтобы соответствующий многочлен p_m принимал как можно меньшие по модулю значения на спектре матрицы B* . То есть, идеальным вариантом был бы характеристический многочлен

$$p_m(t) = (t - \lambda_1)(t - \lambda_2) \dots (t - \lambda_n),$$

но он нам, во-первых, неизвестен, а во-вторых степень его слишком велика, чтобы использовать на практике.

Вместо этого мы вспомним, что $\sigma(B) \subset [-\rho, \rho]$, и попробуем найти «универсальный» многочлен p_m , который бы на всем отрезке $[-\rho, \rho]$ как можно меньше отклонялся от нуля. При этом в силу (2.44) этот многочлен должен удовлетворять условию $p_m(1) = 1$. Говоря более конкретно, необходимо

найти многочлен p_m степени m такой, что $p_m(1) = 1$, а величина

$$\max_{t \in [-\rho, \rho]} |p_m(t)|$$

принимает как можно меньшее значение.

Сразу дадим готовое решение этой задачи:

$$p_m(t) = \frac{T_m(t/\rho)}{T_m(1/\rho)}, \quad (2.46)$$

где T_m — знаменитые *многочлены Чебышева*, о которых стоит рассказать отдельно в небольшом отступлении.

2.8.2. Многочлены Чебышева

Рассмотрим на отрезке $[-1, 1]$ семейство функций $\{T_m\}$, определяемое формулой

$$T_m(t) = \cos(m \arccos t), \quad m = 0, 1, 2, \dots \quad (2.47)$$

Как это ни кажется странным на первый взгляд, $T_m(t)$ при $t \in [-1, 1]$ является алгебраическим многочленом степени m :

$$T_0(t) = 1, \quad T_1(t) = t,$$

а для произвольного m справедливо рекуррентное соотношение

$$T_m(t) = 2t T_{m-1}(t) - T_{m-2}(t), \quad (2.48)$$

которое позволяет вычислять $T_m(t)$ при любых $t \in \mathbb{R}$.

▷₂ Вычислите $T_m(t)$ для $m = 2, 3, 4, 5$.

▷₃ Из (2.47) вывести рекуррентное соотношение (2.48), используя тригонометрическую формулу $\cos \alpha \cos \beta = \frac{1}{2}(\cos(\alpha + \beta) + \cos(\alpha - \beta))$.

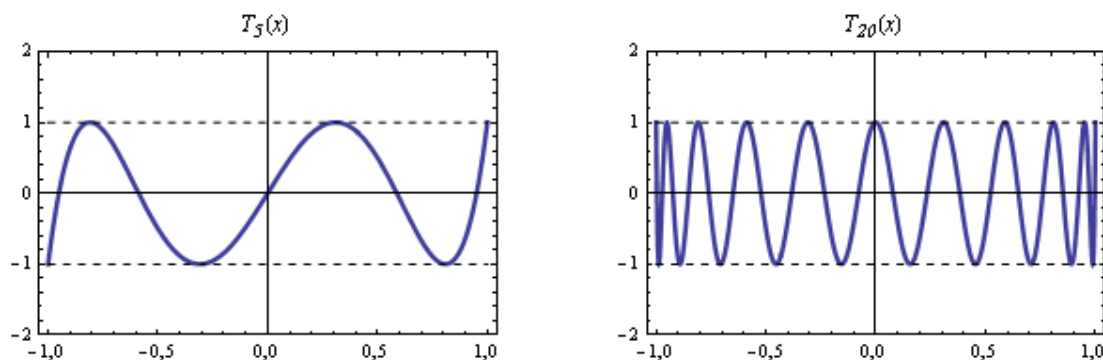


Рис. 2.1. Многочлены Чебышева T_5 и T_{20}

Из рекуррентного соотношения видно, что старший коэффициент многочлена T_m равен 2^{m-1} (при $m > 0$). Основное же свойство многочленов Чебышева состоит в том, что их графики «равномерно колеблются» между -1 и 1 (рис. 2.1). Говоря точнее, многочлен T_m имеет $m + 1$ локальных экстремумов, которые поочередно равны ± 1 .

Следствием этого свойства является то, что *многочлен Чебышева T_m имеет наименьшее отклонение от нуля на отрезке $[-1, 1]$ среди всех многочленов степени m со старшим коэффициентом 2^{m-1}* . Тогда многочлен $T_m(t/\rho)$ будет наименее отклоняться от нуля на $[-\rho, \rho]$ среди всех многочленов, имеющих такой же как у него старший коэффициент. Добавляя нормировочный множитель, в итоге для решения задачи поиска оптимального многочлена p_m получаем формулу (2.46) (см. также рис. 2.2).

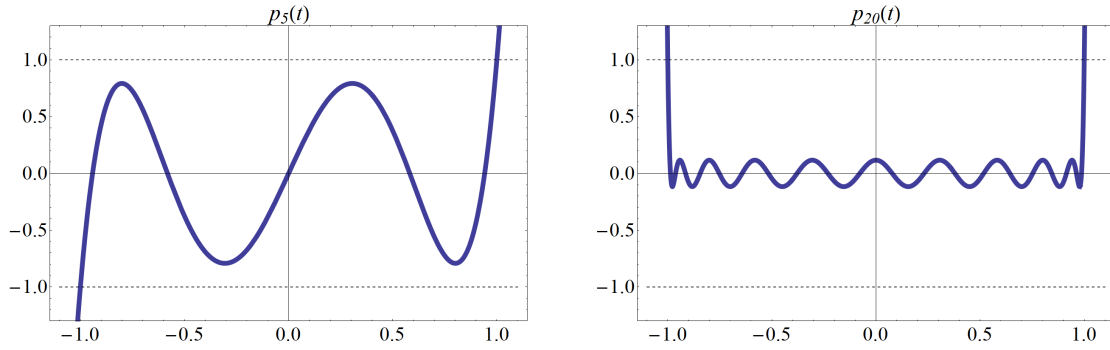


Рис. 2.2. Графики многочленов p_5 и p_{20} для $\rho = 0.99$.

2.8.3. Вывод расчетных формул

Подведем локальный итог вышесказанному: на основе сходящегося стационарного процесса $x^{k+1} = Bx^k + g$ мы строим приближения вида $y^m = \sum_{k=0}^m c_{m,k} x^k$, где $c_{m,k}$ определены согласно (2.46) как коэффициенты многочлена

$$p_m(t) = \sum_{k=0}^m c_{m,k} t^k = \mu_m T_m(t/\rho),$$

$$\mu_m = \frac{1}{T_m(1/\rho)}.$$

Теоретически этой информации достаточно, чтобы для каждого фиксированного m вычислить коэффициенты $c_{m,k}$ и соответствующие значения

y^m . Однако существует гораздо более элегантный и эффективный способ последовательного вычисления y^m , основанный на рекуррентном соотношении (2.48).

Для начала получим на основе этой формулы аналогичные соотношения для многочленов p_m :

$$\begin{aligned} p_m(t) &= \mu_m T_m(t/\rho) = \frac{2\mu_m}{\rho} t T_{m-1}(t/\rho) - \mu_m T_{m-2}(t/\rho) = \\ &= \frac{2\mu_m}{\rho\mu_{m-1}} t p_{m-1}(t) - \frac{\mu_m}{\mu_{m-2}} p_{m-2}(t). \end{aligned} \quad (2.49)$$

Также нам понадобятся рекуррентные соотношения между коэффициентами μ_m :

$$\mu_m^{-1} = T_m(1/\rho) = (2/\rho)T_{m-1}(1/\rho) - T_{m-2}(1/\rho),$$

или

$$\mu_m^{-1} = 2\mu_{m-1}^{-1}/\rho - \mu_{m-2}^{-1}. \quad (2.50)$$

Согласно (2.45) имеем

$$y^m - x^* = p_m(B)(x^0 - x^*),$$

откуда, обозначив

$$\alpha_m = \frac{2\mu_m}{\rho\mu_{m-1}}, \quad \beta_m = -\frac{\mu_m}{\mu_{m-2}}, \quad (2.51)$$

получим

$$\begin{aligned} y^m - x^* &= p_m(B)(x^0 - x^*) = \\ &= (\alpha_m B p_{m-1}(B) + \beta_m p_{m-2}(B))(x^0 - x^*) = \\ &= \alpha_m B (y^{m-1} - x^*) + \beta_m (y^{m-2} - x^*), \end{aligned}$$

или

$$y^m - \alpha_m B y^{m-1} - \beta_m y^{m-2} = (I - \alpha_m B - \beta_m I) x^*. \quad (2.52)$$

Умножая (2.50) на μ_m , получим $\alpha_m + \beta_m = 1$. С учетом тождества $x^* = Bx^* + g$ тогда будем иметь

$$(I - \alpha_m B - \beta_m I) x^* = \alpha_m (I - B) x^* = \alpha_m g.$$

Таким образом (2.52) окончательно принимает чудесный вид

$$y^m = \alpha_m(By^{m-1} + g) + \beta_my^{m-2}, \quad m = 2, 3, \dots \quad (2.53)$$

Это рекуррентное соотношение является основным при программной реализации процесса чебышевского ускорения наряду со стартовым тождеством

$$y^1 = By^0 + g$$

и формулами (2.50), (2.51) для расчета коэффициентов. Если говорить о терминологии, то полученный итерационный процесс можно назвать *двухшаговым нестационарным*.

▷₄ Коэффициенты α_m и β_m можно вычислять по прямым рекуррентным соотношениям без использования промежуточных коэффициентов μ_m . Постройте эти соотношения.

▷₅ Запишите соответствующий вычислительный алгоритм для чебышевских итераций.

2.9. Методы крыловских подпространств

2.9.1. Обобщенный метод минимальных невязок (GMRES)

Определение 2.11. Рассмотрим матрицу A и вектор b . Подпространством Крылова размерности k называется линейная оболочка векторов $b, Ab, A^2b, \dots, A^{k-1}b$:

$$\mathcal{K}_k(A, b) = \text{span}\{b, Ab, A^2b, \dots, A^{k-1}b\}.$$

Подпространства Крылова являются основой многих современных итерационных методов решения СЛАУ и проблемы собственных значений. Сейчас мы познакомимся с одним из наиболее эффективных итерационных методов решения систем с несимметричной матрицей — обобщенным методом минимальных невязок (generalized minimal residual, GMRES).

Суть обобщенного метода минимальных невязок для решения системы $Ax = b$ состоит в следующем: на k -й итерации метода в подпространстве $\mathcal{K}_k(A, b)$ ищется приближенное решение x^k , имеющее минимальную норму невязки $\|Ax^k - b\|_2$. Чтобы сделать формулы более компактными в дальнейшем будем обозначать $x^k = x$.

Так как $x \in \mathcal{K}_k(A, b)$, то его можно представить в виде

$$x = c_1b + c_2Ab + \dots + c_kA^{k-1}b = K_k c,$$

где

$$K_k = [b \mid Ab \mid A^2b \mid \dots \mid A^{k-1}b],$$

а $c = (c_1, \dots, c_k)^T$ — вектор неизвестных коэффициентов, который по построению метода находится путем решения задачи наименьших квадратов

$$\|AK_k c - b\|_2 \rightarrow \min. \quad (2.54)$$

На этом можно уже остановиться, так как задачу наименьших квадратов мы решать уже умеем. Но результат получится не очень хорошим, главным образом из-за плохой обусловленности матрицы K_k . Дальнейшие наши усилия будут направлены на получения качественного и эффективного способа решения задачи (2.54).

Основная идея заключается в том, чтобы вместо столбцов матрицы K_k для разложения искомого вектора x использовать *ортонормированный базис* пространства Крылова $\mathcal{K}_k(A, b)$. Векторы этого базиса будем обозначать q_j , а матрицу, составленную из этих столбцов обозначим Q_k . Тогда вместо $x^k = K_k c$ будем иметь $x^k = Q_k y$, где y — вектор коэффициентов разложения по базису $\{q_j\}$, и задачу минимизации

$$\|AQ_k y - b\|_2 \rightarrow \min.$$

Для построения нужного нам базиса необходимо сначала изучить несколько новых (и вспомнить несколько старых) понятий.

Алгоритм Арнольди

Определение 2.12. Говорят, что квадратная матрица H имеет *форму Хессенберга* (является хессенберговой), если

$$h_{ij} = 0 \quad \forall i \geq j + 2.$$

Другими словами, форма Хессенберга отличается от верхнетреугольной матрицы лишь ненулевыми элементами под главной диагональю:

$$\begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ & \times & \times & \times & \times \\ & & \times & \times & \times \\ & & & \times & \times \end{bmatrix}.$$

Напомним, что преобразование подобия называется преобразование $A \mapsto X^{-1}AX$, где X — произвольная невырожденная матрица. Такие преобразования не изменяют спектр матрицы, так как если $X^{-1}AX\xi = \lambda\xi$, то $A\eta = \lambda\eta$, где $\eta = X\xi$.

Важный факт: любую квадратную вещественную матрицу A можно ортогональными преобразованиями подобия привести к форме Хессенберга, то есть существует такая матрица Q , что $Q^{-1} = Q^T$ и

$$Q^T A Q = H, \quad \text{или} \quad A Q = Q H.$$

Оказывается, что первые k столбцов матрицы Q образуют базис (ортонормированный, естественно) в подпространстве Крылова $\mathcal{K}_k(A, q_1)$.

Лемма 2.4. Пусть $h_{j+1,j} \neq 0$ для $j = 1, \dots, k-1$. Тогда

$$\text{span}\{q_1, \dots, q_k\} = \mathcal{K}_k(A, q_1).$$

Доказательство. При $k = 1$ утверждение леммы очевидно. Предположим теперь, что $\text{span}\{q_1, \dots, q_{k-1}\} = \mathcal{K}_{k-1}(A, q_1)$. В тождестве $AQ = QH$ рассмотрим $(k-1)$ -й столбец матриц AQ и QH с учетом, что H хессенбергова:

$$Aq_{k-1} = h_{1,k-1}q_1 + h_{2,k-1}q_2 + \dots + h_{k,k-1}q_k, \quad (2.55)$$

откуда

$$q_k = \frac{1}{h_{k,k-1}} \left(Aq_{k-1} - \sum_{j=1}^{k-1} h_{j,k-1}q_j \right). \quad (2.56)$$

Применяя рекурсивно эту формулу к q_{k-1} и так далее, легко убедиться, что вектор q_k представим в виде

$$q_k = \gamma_k A^{k-1} q_1 + \delta_k,$$

где $\gamma_k \neq 0$, а $\delta_k \in \text{span}\{q_1, \dots, q_{k-1}\} = \mathcal{K}_{k-1}(A, q_1)$. Следовательно,

$$\text{span}\{q_1, q_2, \dots, q_{k-1}, q_k\} = \text{span}\{q_1, q_2, \dots, q_{k-1}, A^{k-1}q_1\} = \mathcal{K}_k(A, q_1). \quad \square$$

Итак, ортогональный базис подпространства Крылова может быть получен путем нахождения первых k столбцов матрицы Q , которая преобразует матрицу A к форме Хессенберга. Доказательство леммы 2.4 дает и способ нахождения искомых векторов q_j , называемый *алгоритмом Арнольди*.

Опишем один шаг алгоритма. Предположим, что уже найдены векторы q_1, \dots, q_j . Перепишем (2.55), заменив k на $j+1$:

$$Aq_j = h_{1j}q_1 + h_{2j}q_2 + \dots + h_{j+1,j}q_{j+1}.$$

По построению векторы q_i ортонормированы, поэтому все коэффициенты h_{ij} кроме последнего, $h_{j+1,j}$, можно найти, умножив последнее тождество скалярно на q_i :

$$h_{ij} = (Aq_j, q_i), \quad i = 1, \dots, j. \quad (2.57)$$

Искомый вектор q_{j+1} находим аналогично (2.56):

$$q_{j+1} = \frac{1}{h_{j+1,j}} \left(Aq_j - \sum_{i=1}^j h_{ij}q_i \right), \quad (2.58)$$

где в силу того, что $\|q_{j+1}\|_2 = 1$,

$$h_{j+1,j} = \left\| Aq_j - \sum_{i=1}^j h_{ij}q_i \right\|_2. \quad (2.59)$$

Заметим, что коэффициенты h_{ij} представляют собой элементы j -го столбца матрицы H . Таким образом алгоритм Арнольди позволяет последовательно находить столбцы матриц Q и H в разложении $Q^T A Q = H$.

Алгоритм Арнольди	
1)	<i>Задать вектор q_1 такой, что $\ q_1\ _2 = 1$</i>
2)	for $j = 1$ to $k - 1$
3)	· $z = Aq_j$
4)	· for $i = 1$ to j
5)	· · $h_{ij} = (z, q_i)$
6)	· · $z = z - h_{ij}q_i$
7)	· $h_{j+1,j} = \ z\ _2$
8)	· if $h_{j+1,j} = 0$ then stop
9)	· $q_{j+1} = z/h_{j+1,j}$

Необходимо сделать несколько важных комментариев касательно данного алгоритма.

- С формальной точки зрения строка 5 не точно соответствует формуле (2.57), так как при $i > 1$ переменная z уже не будет равна Aq_j . Тем не менее, алгоритм будет давать нужный результат в силу ортогональности векторов q_1, \dots, q_j .

- Обычно необходимое количество векторов-столбцов матрицы Q , которое мы обозначили k , заранее не известно. Условие выхода из внешнего цикла зависит от приложения. В обобщенном методе минимальных невязок, например, это будет условие малости невязки $Ax^k - b$ на текущем приближении.

- В общем случае вектор q_1 может быть произвольным вектором единичной длины. В обобщенном методе минимальных невязок же по построению $q_1 = b/\|b\|_2$.

- В строке 8 осуществляется проверка на «вырожденный» случай, когда невозможно вычислить q_{j+1} и алгоритм останавливается досрочно. Тогда будет справедливо тождество

$$Aq_j = h_{1j}q_1 + \dots + h_{jj}q_j,$$

или $Aq_j \in \mathcal{K}_j(A, q_1)$, откуда следует, что

$$Ax \in \mathcal{K}_j(A, q_1) \quad \forall x \in \mathcal{K}_j(A, q_1),$$

то есть подпространство Крылова размерности j является *инвариантным* относительно матрицы (отображения) A . Как следует из приведенной ниже леммы это означает, что метод GMRES в этом случае даст точное решение.

Лемма 2.5. *Пусть подпространство Крылова $\mathcal{K}_j = \mathcal{K}_j(A, b)$ инвариантно относительно A . Тогда если*

$$x^j = \operatorname{argmin}_{x \in \mathcal{K}_j} \|Ax - b\|_2, \quad (2.60)$$

то $Ax^j = b$.

Доказательство. В силу инвариантности имеем $Ax^j \in \mathcal{K}_j$, при этом $b \in \mathcal{K}_j$ по построению. Значит, задача наименьших квадратов (2.60) представляет собой поиск проекции вектора b на подпространство, которому он и так принадлежит. Поэтому эта проекция, равная Ax^j , совпадает с b . \square

Таким образом, «вырожденный» случай в алгоритме Арнольди является счастливым для обобщенного метода минимальных невязок.

Вычислительная схема метода GMRES

Теперь нужно собрать вместе все детали. Пусть известна матрица $Q_{k-1} = [q_1 | \dots | q_k]$, столбцы которой образуют ортонормированный базис в подпространстве $\mathcal{K}_k(A, b)$, при этом $q_1 = b/\|b\|_2$. Тогда для вычисления приближения x^k необходимо:

1) Выполнить один шаг алгоритма Арнольди для вычисления коэффициентов h_{ik} (строки 3-9 при $j = k$). Эти коэффициенты нужны для решения задачи наименьших квадратов в следующем пункте.

2) Найти $y^k = \operatorname{argmin}_y \|AQ_k y - b\|_2$ (ниже будет сказано о том, как сделать это эффективно).

3) Вычислить $x^k = Q_k y^k$.

4) Если норма невязки $\|Ax^k - b\|_2$ превышает требуемое значение, положить $k = k + 1$ и перейти к шагу 1.

Задача наименьших квадратов в пункте 2 имеет размерность $n \times k$, однако ее можно свести к эквивалентной задаче, размерность которой будет всего $(k + 1) \times k$. Покажем как это делается.

В силу ортогональности матрицы Q имеем

$$\|AQ_k y - b\|_2 = \|Q^T AQ_k y - Q^T b\|_2.$$

Так как q_1 представляет собой нормированный вектор b , и все строки Q^T взаимно ортогональны, то $Q^T b = d$,

$$d = (\|b\|_2, 0, \dots, 0)^T. \quad (2.61)$$

Далее, по построению имеем $Q^T A Q = H$. Напомним, что ненулевые элементы столбцов матрицы H находятся в процессе выполнения алгоритма Арнольди. Нетрудно заметить, что матрица $Q^T A Q_k$ состоит из первых k столбцов матрицы H , причем строки с номерами от $k + 2$ до n будут нулевыми, так как H хессенбергова. А так как в вектор d только первая компонента отлична от нуля, получаем

$$\|Q^T AQ_k y - Q^T b\|_2 = \|Q^T AQ_k y - d\|_2 = \|\tilde{H}_k y - d_{k+1}\|_2,$$

где \tilde{H}_k — матрица, составленная из первых k столбцов и $(k + 1)$ строк матрицы H , а $d_{k+1} = (\|b\|_2, 0, \dots, 0)^T \in \mathbb{R}^{k+1}$. Таким образом, на шаге 2) алгоритма GMRES решается задача наименьших квадратов с хессенберговой матрицей:

$$y^k = \operatorname{argmin}_{y \in \mathbb{R}^k} \|\tilde{H}_k y - d_{k+1}\|_2. \quad (2.62)$$

Для решения этой задачи нужно построить QR -разложение матрицы \tilde{H} , что осуществляется за $O(k^2)$ операций с помощью преобразований вращения (см. далее). Кроме этого, если воспользоваться имеющимся с

предыдущей итерации QR -разложением матрицы \tilde{H}_{k-1} , то объем работы можно уменьшить на порядок, до $O(k)$ операций.

▷₁ Подумайте как это сделать.

В заключение обсуждения этого метода заметим, что по построению при отсутствии ошибок округления метод GMRES сходится не более чем за n итераций (n — размер матрицы A).

2.9.2. Метод полной ортогонализации (FOM)

Метод полной ортогонализации (full orthogonalization method, FOM) очень похож на обобщенный метод минимальных невязок. Приближенное решение x^k так же ищется в подпространстве Крылова и раскладывается по ортонормированному базису $\{q_j\}$, $x^k = Q_k y^k$, но вместо условия минимизации невязки в среднеквадратичном смысле накладывается условие ортогональности

$$(Ax^k - b) \perp \mathcal{K}_k(A, b). \quad (2.63)$$

То есть вектор невязки должен быть ортогонален всему подпространству Крылова, что эквивалентно ортогональности всем векторам из базиса этого подпространства. В качестве такого базиса, конечно, следует брать ортонормированный базис $\{q_j\}$, построенный методом Арнольди. В этом случае условие (2.63) примет вид

$$(AQ_k y^k - b, q_j) = 0, \quad j = \overline{1, k},$$

или

$$Q_k^T A Q_k y^k = Q_k^T b.$$

По построению имеем

$$Q_k^T A Q_k = H_k,$$

где $H_k = [H]_k$ — первые k строк и столбцов матрицы H .

Таким образом, алгоритм метода FOM отличается от метода GMRES только пунктом 2: вместо задачи наименьших квадратов здесь нужно решать СЛАУ

$$H_k y^k = d_k, \quad (2.64)$$

где H_k хессенбергова матрица размера k , а вектор $d_k \in \mathbb{R}^k$ определен аналогично (2.61).

▷₂ Запишите алгоритм решения СЛАУ с матрицей в форме Хессенберга а) методом Гаусса, б) методом вращений.

Глава 3

Методы решения проблемы собственных значений

- 3.1. Проблема собственных значений: общая характеристика
- 3.2. Степенной метод
- 3.3. Метод Данилевского
- 3.4. QR -алгоритм
- 3.5. Метод вращений Якоби

3.1. Проблема собственных значений: общая характеристика

3.1.1. Сведения из линейной алгебры

Еще раз вспомним необходимые сведения из курса линейной алгебры.

Пусть A — вещественная квадратная матрица. Ненулевой вектор v называется *собственным вектором* матрицы A , если существует $\lambda \in \mathbb{C}$ такое, что

$$Av = \lambda v.$$

Число λ называется *собственным значением*, соответствующим v . Множество всех собственных значений A называется *спектром* и обозначается $\sigma(A)$.

Очевидно, что комплексным собственным значениям соответствуют комплексные собственные векторы.

Замечание 3.1. Собственный вектор определен с точностью до постоянного множителя:

$$Av = \lambda v \quad \Rightarrow \quad A(\alpha v) = \lambda(\alpha v).$$

Поэтому все собственные векторы, соответствующие собственному значению λ , образуют так называемое *собственное подпространство*. Размерность собственного подпространства (число линейно независимых собственных векторов с собственным значением λ) называют *геометрической кратностью* собственного значения.

Как известно, все собственные значения являются корнями *характеристического многочлена*

$$P(\lambda) = \det(A - \lambda I).$$

Если λ — корень многочлена P кратности k , то говорят, что *алгебраическая кратность* λ равна k .

Известен следующий факт (теорема Абеля — Руффини): для любой степени $n \geq 5$ существует многочлен P_n с рациональными коэффициентами, вещественный корень которого не может быть представлен в виде выражения, составленного из рациональных чисел и операций сложения, вычитания, умножения, деления и корня целой степени. Таким образом, проблема собственных значений имеет принципиальное отличие от задачи решения СЛАУ. Решение

СЛАУ с рациональными коэффициентами всегда является вектором рациональных чисел и может быть найдено точно при отсутствии ошибок округления, например, тем же методом Гаусса. Собственные же значения матрицы в общем случае не могут быть найдены точно, поэтому *все методы вычисления собственных значений являются итерационными*.

Определение 3.1. Если квадратная матрица размерности n имеет n линейно независимых собственных векторов, то она называется *диагонализируемой*, а соответствующий ей линейный оператор — *оператором простой структуры*.

Эквивалентным определением можно считать следующее: матрица диагонализируема если и только если она может быть приведена к диагональному виду преобразованием подобия:

$$V^{-1}AV = \Lambda, \quad \det V \neq 0, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n).$$

Нетрудно убедиться, что столбцы матрицы V являются собственными векторами матрицы A .

Задачи на нахождение собственных значений условно делятся на два класса. Если нужно найти одно или несколько собственных значений и соответствующие им собственные векторы, то проблема собственных значений называется *частичной*. Если же необходимо найти все собственные значения и векторы, проблема называется *полной*.

3.2. Степенной метод

Степенной метод позволяет найти максимальное по модулю собственное значение и соответствующий ему собственный вектор вещественной диагонализируемой матрицы. Он использовался, например, в алгоритме PageRank при расчете релевантности веб-страницы поисковому запросу.

Пусть A — диагонализируемая матрица, $\lambda_1, \dots, \lambda_n$ — упорядоченные по убыванию модуля собственные значения, v_1, \dots, v_n — соответствующий базис из собственных векторов. В дальнейшем будем предполагать, что евклидова норма всех v_i равна единице. Рассмотрим несколько возможных случаев.

3.2.1. Случай 1

Пусть $\lambda_1 \in \mathbb{R}$, $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$.

Разложим по базису $\{v_i\}$ произвольный вектор $y^0 \in \mathbb{R}^n$: $y^0 = \sum_1^n \alpha_i v_i$. Предположим, что $\alpha_1 \neq 0$ и рассмотрим последовательность (y^k) :

$$y^{k+1} = Ay^k.$$

Тогда

$$y^k = A^k y^0 = \sum_{i=1}^n \alpha_i \lambda_i^k v_i = \lambda_1^k \left(\alpha_1 v_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k v_i \right). \quad (3.1)$$

Значит, при $k \rightarrow \infty$ имеем

$$y^k = \lambda_1^k (\alpha_1 v_1 + O(\lambda_2^k / \lambda_1^k)),$$

т. е. вектор y^k все сильнее приближается к собственному подпространству, соответствующему v_1 . Для практического применения необходимо на каждом шаге нормировать y^k . Как правило, y^k делят на его евклидову норму, то есть фактически вычисляется последовательность

$$u^k = y^k / \|y^k\|_2 = A^k y^0 / \|A^k y^0\|_2. \quad (3.2)$$

Такой подход облегчает вычисление собственного значения λ_1 с помощью *отношения Релея*, для знакомства с которым необходимо сделать небольшое отступление.

Пусть y — некоторое приближение к собственному вектору матрицы A . Тогда, скорее всего, не существует такого числа λ , для которого $Ay = \lambda y$. Но можно определить «почти собственное число» как скаляр λ , для которого норма невязки $Ay - \lambda y$ минимальна:

$$r(y) = \operatorname{argmin}_{\lambda \in \mathbb{C}} \| \lambda y - Ay \|_2.$$

Это обычная одномерная задача наименьших квадратов, решая которую приходим к следующему определению.

Определение 3.2. Для любого вектора $y \neq 0$ отношением Релея называется число

$$r(\lambda) = \frac{(y, Ay)}{(y, y)}. \quad (3.3)$$

▷₁ Выведите эту формулу.

Исходя из вышесказанного легко увидеть, что если y является собственным вектором, то $r(y)$ будет в точности равно его собственному значению.

Можно показать, что если $\|y - v\| \leq \varepsilon$, то $|r(y) - \lambda| = O(\varepsilon)$, а в случае $A = A^T$ будем иметь квадратичную сходимость: $|r(y) - \lambda| = O(\varepsilon^2)$.

▷₂ Докажите эти утверждения.

Таким образом, получаем следующий **базовый алгоритм степенного метода**:

```
 $u \leftarrow y^0 / \|y^0\|_2; \quad \lambda \leftarrow (u, Au)$   
while  $\|Au - \lambda u\| > \varepsilon$  do  
     $y \leftarrow Au$   
     $u \leftarrow y / \|y\|_2$   
     $\lambda \leftarrow (u, Au)$   
end while
```

На выходе будем иметь $u \approx \pm v_1$, $\lambda \approx \lambda_1$. Здесь ε — требуемая точность.

▷₃ Объясните, почему мы вынуждены писать $u \approx \pm v_1$? В каких случаях будем иметь $u \approx v_1$, а в каких $u \approx -v_1$?

▷₄ Оцените количество арифметических операций в одной итерации степенного метода.

Из (3.1) видно, что степенной метод сходится со скоростью геометрической прогрессии со знаменателем $|\lambda_2/\lambda_1|$.

▷₅ Докажите это утверждение строго.

3.2.2. Случай 2

Пусть $\lambda_1 = \lambda_2 = \dots = \lambda_m \in \mathbb{R}$, $|\lambda_1| > |\lambda_{m+1}| \geq \dots \geq |\lambda_n|$.

Так как матрица A по условию диагонализируема, геометрическая кратность λ_1 равна m , т. е. собственные векторы v_1, \dots, v_m линейно независимы и образуют собственное подпространство V_m размерности m . Тогда формула (3.1) примет вид

$$y^k = \lambda_1^k \left(\sum_{i=1}^m \alpha_i v_i + \sum_{i=m+1}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k v_i \right) \approx \lambda_1^k \sum_{i=1}^m \alpha_i v_i.$$

В этом случае алгоритм остается без изменений. Последовательность (u^k) сходится к какому-то вектору $u \in V_m$ — он ничем не хуже других собственных векторов v_1, \dots, v_m .

3.2.3. Случай 3

Пусть $\lambda_1 = -\lambda_2 \in \mathbb{R}$, $|\lambda_1| > |\lambda_3| \geq \dots \geq |\lambda_n|$. Не нарушая общности можно считать, что $\lambda_1 > 0$.

В этом случае (3.1) превращается в

$$\begin{aligned} y^k &= \sum_{i=1}^n \alpha_i \lambda_i^k v_i = \alpha_1 \lambda_1^k v_1 + \alpha_2 (-\lambda_1)^k v_2 + \sum_{i=3}^n \alpha_i \lambda_i^k v_i = \\ &= \lambda_1^k \left(\alpha_1 v_1 + (-1)^k \alpha_2 v_2 + \sum_{i=3}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k v_i \right), \end{aligned}$$

откуда получим

$$\begin{aligned} y^{2k} &\approx \lambda_1^{2k} (\alpha_1 v_1 + \alpha_2 v_2), \\ y^{2k+1} &\approx \lambda_1^{2k+1} (\alpha_1 v_1 - \alpha_2 v_2), \end{aligned} \tag{3.4}$$

т. е. последовательность (u^k) (3.2) не сходится, а распадается на две сходящиеся подпоследовательности (u^{2k}) и (u^{2k+1}) . Эти последовательности сходятся к двум различным векторам из $\text{span}(v_1, v_2)$, причем в отличие от предыдущего случая, ни один из них не является, вообще говоря, собственным.

Но выход есть. Рассмотрим подпоследовательность четных элементов (u^{2k}) ,

$$u^{2k} = y^{2k} / \|y^{2k}\|_2, \quad k = 1, 2, \dots$$

После достаточно большого числа итераций будем иметь

$$u^{2k} \approx \beta_1 v_1 + \beta_2 v_2, \quad \|u^{2k}\|_2 = 1.$$

Тогда

$$A^2 u^{2k} \approx \lambda_1^2 \beta_1 v_1 + \lambda_1^2 \beta_2 v_2 = \lambda_1^2 u^{2k},$$

то есть

$$(u^{2k}, A^2 u^{2k}) \rightarrow \lambda_1^2. \tag{3.5}$$

Таким образом находим λ_1 .

Теперь рассмотрим, как можно найти соответствующий собственный вектор. Имеем

$$\begin{aligned} u^{2k} &\approx \beta_1 v_1 + \beta_2 v_2, \\ Au^{2k} &\approx \lambda_1 \beta_1 v_1 - \lambda_1 \beta_2 v_2, \\ A^2 u^{2k} &\approx \lambda_1^2 \beta_1 v_1 + \lambda_1^2 \beta_2 v_2 \end{aligned}$$

Так как приближенное значение λ_1 уже известно, можно исключить компоненту v_2 , складывая последнее равенство с предыдущим, умноженным на λ_1 :

$$\lambda_1 Au^{2k} + A^2 u^{2k} \approx 2\lambda_1^2 \beta_1 v_1.$$

При необходимости пронормировав это выражение, получим приближение к v_1 . Аналогично будем иметь

$$\lambda_1 Au^{2k} - A^2 u^{2k} \approx -2\lambda_1^2 \beta_1 v_2.$$

▷₆ Запишите соответствующий алгоритм приближенного вычисления λ_1, v_1, v_2 .

3.2.4. Комплексные собственные значения

Покажем теперь как можно определить сразу два максимальных по модулю собственных значения, в том числе и комплексных. Пусть $|\lambda_1| \geq |\lambda_2| > |\lambda_3| > \dots > |\lambda_n|$ и в разложении начального приближения y^0 по собственным векторам коэффициенты при v_1 и v_2 отличны от 0. Тогда после достаточно большого числа итераций будем иметь

$$\begin{aligned} u^k &\approx \beta_1 v_1 + \beta_2 v_2 = \eta_0, \\ Au^k &\approx \lambda_1 \beta_1 v_1 + \lambda_2 \beta_2 v_2 = \eta_1, \\ A^2 u^k &\approx \lambda_1^2 \beta_1 v_1 + \lambda_2^2 \beta_2 v_2 = \eta_2. \end{aligned}$$

Возьмем по вектору с каждой стороны этих приближенных равенств и составим из них матрицы:

$$K = [u^k | Au^k | A^2 u^k] \approx [\eta_0 | \eta_1 | \eta_2]. \quad (3.6)$$

Векторы η_i принадлежат двумерному подпространству $\text{span}\{v_1, v_2\}$ и, следовательно, линейно зависимы. Это означает, что существует вектор $c = (c_0, c_1, c_2)^T \neq 0$ такой, что

$$c_0 \eta_0 + c_1 \eta_1 + c_2 \eta_2 = [\eta_0 | \eta_1 | \eta_2] c = 0.$$

Учитывая определение η_i , запишем это тождество в виде

$$\begin{aligned} [\eta_0 | \eta_1 | \eta_2] c &= \underbrace{[\beta_1 v_1 | \beta_2 v_2]}_B \begin{bmatrix} 1 & \lambda_1 & \lambda_1^2 \\ 1 & \lambda_2 & \lambda_2^2 \end{bmatrix} c = \\ &= B \begin{bmatrix} c_0 + c_1 \lambda_1 + c_2 \lambda_1^2 \\ c_0 + c_1 \lambda_2 + c_2 \lambda_2^2 \end{bmatrix} = B \begin{bmatrix} P(\lambda_1) \\ P(\lambda_2) \end{bmatrix} = 0, \end{aligned}$$

где

$$P(t) = c_2 t^2 + c_1 t + c_0.$$

Так как столбцы матрицы B линейно независимы, отсюда следует, что

$$P(\lambda_1) = P(\lambda_2) = 0,$$

то есть искомые собственные значения можно найти, решив квадратное уравнение $P(\lambda) = 0$. Осталось найти коэффициенты многочлена P . Они определены с точностью до постоянного множителя, поэтому можно положить $c_2 = 1$. Остальные коэффициенты найдем, используя исходное тождество (3.6):

$$Kc = c_0 u^k + c_1 A u^k + A^2 u^k \approx [\eta_0 | \eta_1 | \eta_2] c = 0,$$

Отсюда можно приближенно найти коэффициенты многочлена P путем решения задачи наименьших квадратов

$$(c_0, c_1)^T = \operatorname{argmin}_{c \in \mathbb{R}^2} \left\| [u^k | A u^k] c + A^2 u^k \right\|_2. \quad (3.7)$$

После этого, решая квадратное уравнение $P(t) = 0$, находятся приближения к искомым собственным значениям λ_1 и λ_2 .

▷₇ Как теперь можно найти приближения к собственным векторам v_1 и v_2 ?

3.2.5. Обратная итерация со сдвигом

Идею степенного метода можно использовать для определения любого, а не только максимального по модулю собственного значения. Рассмотрим скаляр $\mu \notin \sigma(A)$. Тогда матрица $(A - \mu I)^{-1}$ имеет такие же собственные векторы, как и A , а соответствующие собственные значения ее равны

$$\mu_i = (\lambda_i - \mu)^{-1}.$$

▷₈ Докажите это.

Отсюда видно, что, применив к матрице $(A - \mu I)^{-1}$ степенной метод, мы можем найти ближайшее к μ собственное значение λ_i , а также соответствующий ему собственный вектор. При этом нет необходимости обращать матрицу $A - \mu I$: для вычисления $y = (A - \mu I)^{-1}u^k$ на каждой итерации нужно лишь решать СЛАУ

$$(A - \mu I)y = u^k.$$

Для снижения вычислительных затрат можно предварительно построить LU - или QR - разложение матрицы $A - \mu I$. Общий алгоритм метода обратной итерации имеет вид

```
u ← y0 / ||y0||2
λ ← (u, Au)
while ||Au - λu|| > ε do
    Решить систему (A - μI)y = u
    u ← y / ||y||2
    λ ← (u, Au)
end while
```

Обратная итерация с отношением Релея

Алгоритм обратной итерации естественным образом можно улучшить, исходя из следующего соображения: так как скорость сходимости метода зависит от близости параметра μ к какому-то собственному значению λ_j , можно на каждой итерации уточнять μ , заменяя его на текущее приближение к λ_j . Приведенный ниже алгоритм лишь на один символ отличается от предыдущего, но имеет на порядок более высокую скорость сходимости.

```
u ← y0 / ||y0||2
λ ← (u, Au)
while ||Au - λu|| > ε do
    Решить систему (A - λI)y = u
    u ← y / ||y||2
    λ ← (u, Au)
end while
```

3.3. Метод Данилевского

3.3.1. Преобразования подобия

Многие методы решения проблемы собственных значений используют принцип, аналогичный прямым методам решения СЛАУ: исходная матрица A путем некоторых преобразований приводится к эквивалентной матрице A' , для которой проблема собственных значений решается просто. Естественно, решения задачи для матриц A и A' должны совпадать, т. е. преобразование $A \mapsto A'$ должно сохранять спектр. Другими словами, это преобразование должно быть *преобразованием подобия*.

Определение 3.3. Пусть S — невырожденная матрица. Преобразованием подобия квадратной матрицы A называется преобразование

$$A \mapsto S^{-1}AS.$$

Из курса линейной алгебры известно, что преобразование подобия сохраняет характеристический многочлен, а следовательно, и спектр матрицы. Следует понимать, однако, что такое преобразование изменяет собственные векторы:

$$A'x' = \lambda x' \Leftrightarrow S^{-1}ASx' = \lambda x' \Leftrightarrow ASx' = \lambda Sx',$$

т. е. если x' — собственный вектор A' , то ему соответствует собственный вектор $x = Sx'$ матрицы A .

Таким образом, можно сформулировать следующее **общее правило**: если к матрице A применяется преобразование вида

$$A \mapsto TA = A',$$

то для сохранения спектра необходимо после этого выполнить преобразование

$$A' \mapsto A'T^{-1},$$

и наоборот.

В качестве преобразования T , как правило, используются те же преобразования, что и для СЛАУ, т. е., в частности, элементарные преобразования. Напомним, что эти преобразования бывают двух типов. Элементарное преобразование *первого рода* — это умножение i -й строки матрицы A на произвольный скаляр $\alpha \neq 0$. Данное преобразование

эквивалентно умножению слева на *элементарную матрицу первого типа*

$$T_i(\alpha) = \text{diag}(1, \dots, 1, \alpha, 1, \dots, 1),$$

где элемент α находится на i -й позиции.

Элементарное преобразование *второго рода* — это добавление к i -й строке матрицы A ее j -й строки, умноженной на α . Такое преобразование задается матрицей $S_{ij}(\alpha)$, которая представляет собой матрицу с единицами на главной диагонали и элементом α на позиции (i, j) . Эту матрицу будем называть *элементарной матрицей второго типа*.

Чтобы в дальнейшем можно было использовать эти преобразования для решения проблемы собственных значений, нам необходимо (в соответствии с приведенным выше общим правилом) разобраться с тем, как действуют соответствующие операции «дополнения до подобия», т. е. чему эквивалентны операции

$$A \mapsto A(T_i(\alpha))^{-1} \quad \text{и} \quad A \mapsto A(S_{ij}(\alpha))^{-1}.$$

Элементарное преобразование подобия первого типа. Для начала вычислим обратную к матрице $T(\alpha_i)$. Очевидно, что

$$(T_i(\alpha))^{-1} = T_i(\alpha^{-1}).$$

Также нетрудно убедиться, что преобразование $A \mapsto AT_i(\alpha^{-1})$ представляет собой умножение i -го столбца матрицы A на α^{-1} . Таким образом, *если i -я строка (столбец) матрицы A умножается на скаляр α , то для сохранения спектра необходимо после этого разделить i -й столбец (строку) полученной матрицы на α* . Описанное преобразование будем называть *элементарным преобразованием подобия первого типа*.

Элементарное преобразование подобия второго типа. Обратная матрица к $S_{ij}(\alpha)$ вычисляется легко:

$$(S_{ij}(\alpha))^{-1} = S_{ij}(-\alpha).$$

Выяснить, что происходит с матрицей A при умножении справа на $S_{ij}(-\alpha)$, нам поможет известное тождество

$$(AB)^T = B^T A^T.$$

Имеем

$$AS_{ij}(-\alpha) = (S_{ij}(-\alpha)^T A^T)^T = (S_{ji}(-\alpha) A^T)^T,$$

т.е. если к i -й строке (столбцу) матрицы A прибавлена j -я ее строка (столбец), умноженная на α , то для сохранения спектра необходимо после этого вычесть из j -го столбца (строки) полученной матрицы i -й столбец (строку), умноженный на α .

Преобразование подобия с матрицей перестановки. Помимо элементарных преобразований мы также будем использовать преобразование перестановки местами строк i и j , которое задается матрицей P_{ij} . Эта матрица представляет собой единичную матрицу с переставленными i -й и j -й строками. Легко видеть, что

$$P_{ij}^{-1} = P_{ij}.$$

Кроме этого,

$$AP_{ij} = (P_{ij}A^T)^T,$$

поэтому третье правило преобразований подобия звучит так: *если в матрице переставлены строки (столбцы) с номерами i и j , то для сохранения спектра необходимо после этого у полученной матрицы переставить i -й и j -й столбцы (строки).*

Теперь мы готовы приступить к изучению метода Данилевского.

3.3.2. Метод Данилевского

Метод Данилевского позволяет вычислить **характеристический многочлен**

$$P(\lambda) = |A - \lambda I|$$

для произвольной квадратной матрицы A .

Рассмотрим матрицу вида

$$\begin{bmatrix} p_1 & p_2 & \dots & p_{n-1} & p_n \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix}. \quad (3.8)$$

Такой вид матрицы называется *формой Фробениуса*. Ее характеристический многочлен легко вычисляется путем рекурсивного разложения

определителя по столбцу:

$$\begin{vmatrix} p_1 - \lambda & p_2 & \dots & p_{n-1} & p_n \\ 1 & -\lambda & 0 & \dots & 0 \\ 0 & 1 & -\lambda & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & -\lambda \end{vmatrix} = (p_1 - \lambda)(-\lambda)^{n-1} - \begin{vmatrix} p_2 & \dots & p_{n-1} & p_n \\ 1 & -\lambda & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & 1 & -\lambda \end{vmatrix} = \\ = \dots = (-1)^n (\lambda^n - p_1 \lambda^{n-1} - \dots - p_{n-1} \lambda - p_n). \quad (3.9)$$

Идея метода Данилевского проста: с помощью элементарных преобразований подобия привести данную матрицу A к форме Фробениуса A' . Поскольку преобразования подобия сохраняют спектр матрицы, характеристические многочлены матриц A и A' будут совпадать, т. е. искомый характеристический многочлен может быть вычислен по формуле (3.9).

Рассмотрим алгоритм на примере матрицы A размерности 4. Будем последовательно приводить строки матрицы к нужному виду, начиная с последней.

1. Для начала «делаем единицу» на позиции (4,3): делим третий столбец на $\alpha = a_{43}$. Чтобы сохранить спектр в соответствии с полученным выше правилом, нужно дополнить это преобразование до преобразования подобия, т. е. умножить третью строку на α . В результате указанных операций получим матрицу

$$\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & 1 & \times \end{bmatrix}.$$

2. «Делаем нули» в последней строке: для $j \neq 3$ вычтем из j -го столбца 3-й, умноженный на $\alpha = a_{4j}$. Каждая такая операция должна быть дополнена до преобразования подобия, для чего необходимо к 3-й строке добавлять j -ю, умноженную на α . Получим в итоге

$$\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

3. Аналогично поступим с третьей и второй строками. При этом сделанные ранее нули и единицы не «портятся»:

$$\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & 0 & 1 & 0 \end{bmatrix} \mapsto \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \mapsto \begin{bmatrix} \times & \times & \times & \times \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} p_1 & p_2 & p_3 & p_4 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

Базовый алгоритм метода Данилевского

```

1: for  $i = \overline{n, 2}$  do
2:    $c \leftarrow a_{i,i-1}$ 
3:    $a_{i-1} \leftarrow a_{i-1}/c$ ;
4:    $\underline{a_{i-1}} \leftarrow c \underline{a_{i-1}}$ 
5:   for  $j \neq i$  do
6:      $c \leftarrow a_{ij}$ 
7:      $a_j \leftarrow a_j - ca_{i-1}$ 
8:      $\underline{a_{i-1}} \leftarrow \underline{a_{i-1}} + c \underline{a_j}$ 
9:   end for
10: end for

```

Напомним, что здесь a_j — j -й вектор-столбец матрицы A ; $\underline{a_i}$ — i -я строка матрицы A ; n — размерность матрицы A .

Для эффективной программной реализации строк 6 и 7, как и в методе Гаусса, следует учитывать, что с ходом преобразований нижняя часть матрицы будет содержать все больше нулевых элементов. После выполнения приведенного алгоритма (в невырожденном случае) на месте матрицы A будет находиться ее форма Фробениуса.

Выбор главного элемента. Для минимизации погрешностей округления при машинной реализации необходимо выбирать главный элемент. Для этого перед началом i -го шага нужно выбрать максимальный по модулю элемент a_{ij^*} среди a_{ij} для $j = \overline{1, i-1}$. После этого меняем местами столбцы $i-1$ и j^* , а также соответствующие строки (для сохранения спектра).

Вырожденный случай. Предположим, на i -м этапе не удастся выбрать главный элемент. Это означает, что матрица имеет вид

$$A = \left[\begin{array}{c|c} A_1 & \boxtimes \\ \hline 0 & A_2 \end{array} \right],$$

где A_1, A_2 — квадратные блоки размерности $n-i$ и i соответственно,

причем A_2 имеет форму Фробениуса. Тогда имеем

$$|A - \lambda I| = |A_1 - \lambda I_1| \cdot |A_2 - \lambda I_2|.$$

Здесь I_1 и I_2 — единичные матрицы, размерность которых соответствует блокам A_1 и A_2 . Второй множитель мы можем вычислить сразу, поэтому остается лишь привести к форме Фробениуса матрицу A_1 .

3.4. QR -алгоритм

3.4.1. Общая схема QR -алгоритма

QR -алгоритм является одним из наиболее известных методов для вычисления всех (в том числе и комплексных) собственных значений произвольной квадратной матрицы. Суть алгоритма проста. Построим последовательность матриц

$$A_0 = A, A_1, A_2, \dots,$$

по следующему правилу:

- 1) строится QR -разложение для матрицы A_k : $A_k = Q_k R_k$;
- 2) вычисляется $A_{k+1} = R_k Q_k$.

Описанное преобразование $A_k \mapsto A_{k+1}$ является преобразованием подобия:

$$A_{k+1} = R_k Q_k = Q_k^{-1} A_k Q_k,$$

поэтому все матрицы A_k подобны исходной матрице A .

Теорема 3.1 (QR -алгоритм, упрощенная формулировка). *Последовательность матриц $\{A_k\}_{k=0}^{\infty}$, построенная по описанному выше правилу, при выполнении определенных условий² сходится к почти верхнетреугольной матрице \tilde{R} , у которой на диагонали стоят либо λ_i — собственные значения A , либо блоки вида $\begin{bmatrix} \times & \times \\ \times & \times \end{bmatrix}$, собственные значения которых равны комплексно-сопряженным собственным значениям матрицы A .*

Замечание 3.2. Коэффициенты блоков $\begin{bmatrix} \times & \times \\ \times & \times \end{bmatrix}$ не обязаны сходиться к какому-то пределу, сходятся их собственные значения.

Замечание 3.3. QR -алгоритм *никогда* не применяется к «неподготовленной» матрице A , потому что это слишком трудоемко. Сначала матрицу необходимо преобразованием подобия (желательно, ортогональным) привести к форме Хессенберга: $H = Q^T A Q$ и лишь потом начинать итерации QR -алгоритма. Это снизит трудоемкость одной итерации с $O(n^3)$ до $O(n^2)$.

²Более подробно результаты о сходимости QR -алгоритма обсуждаются в [5, с. 138].

3.4.2. Приведение к форме Хессенберга методом отражений

Нам уже известен один способ ортогонального приведения матрицы к форме Хессенберга — алгоритм Арнольди. Сейчас мы изучим еще один способ, более трудоемкий, но при этом обладающий лучшей вычислительной устойчивостью.

Преобразование матрицы A будем осуществлять по столбцам слева направо. Имеет место почти полная аналогия с процедурой построения QR -разложения из раздела 2.4.3. Различия будет всего два: во-первых, вместо верхнетреугольной теперь нам нужна хессенбергова матрица; во-вторых, нужно использовать преобразования подобия, то есть вместо отображений вида $A \mapsto Q_k A$ будут использоваться отображения $A \mapsto Q_k A Q_k^T$.

Замечание 3.4. Здесь мы используем стандартное обозначение для ортогональных матриц — Q_k . Необходимо понимать, что эти матрицы не имеют никакого отношения к ортогональным матрицам из QR -алгоритма, который был описан в п. 3.4.1.

Итак, на k -м шаге алгоритма в матрице A обнуляются элементы k -го столбца в строках от $k + 2$ до n путем применения преобразования отражения, что соответствует применению отображения

$$A \mapsto Q_k A, \quad Q_k = \begin{bmatrix} I_k & 0 \\ 0 & H_k \end{bmatrix}. \quad (3.10)$$

Здесь I_k — единичная матрица размерности k , H_k — матрица отражения размерности $n - k$, $H_k = I - 2w_k w_k^T$. Преобразование Q_k изменяет в матрице A строки с номерами от $k + 1$ до n . После применения этого преобразования необходимо домножить полученную матрицу справа на Q_k^T , чтобы сохранить спектр:

$$Q_k A \mapsto Q_k A Q_k^T.$$

Это преобразование изменит в матрице $Q_k A$ *столбцы* от $(k + 1)$ -го до n -го, гарантируя сохранение сделанных на предыдущих этапах нулевых элементов.

Рассмотрим для наглядности случай $n = 5$.

$$\begin{aligned}
 k = 1 : & \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ * & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ * & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ * & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ * & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \end{bmatrix} \xrightarrow{Q_1(\cdot)} \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \end{bmatrix} \xrightarrow{(\cdot)Q_1^T} \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \end{bmatrix} ; \\
 k = 2 : & \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & * & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & * & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & * & \mathbf{x} & \mathbf{x} & \mathbf{x} \end{bmatrix} \xrightarrow{Q_2(\cdot)} \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} \end{bmatrix} \xrightarrow{(\cdot)Q_2^T} \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} \end{bmatrix} ; \\
 k = 3 : & \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & * & \mathbf{x} & \mathbf{x} \\ 0 & 0 & * & \mathbf{x} & \mathbf{x} \end{bmatrix} \xrightarrow{Q_3(\cdot)} \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & 0 & \mathbf{x} & \mathbf{x} \end{bmatrix} \xrightarrow{(\cdot)Q_3^T} \begin{bmatrix} \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & \mathbf{x} & \mathbf{x} & \mathbf{x} \\ 0 & 0 & 0 & \mathbf{x} & \mathbf{x} \end{bmatrix} .
 \end{aligned}$$

На этой схеме жирными крестиками обозначены элементы, которые были изменены в матрице в ходе последнего преобразования. Звездочки на k -м шаге обозначают элементы вектора a_k , который служит для определения преобразования отражения H_k из (3.10) по уже знакомой формуле (2.28):

$$w_k = (a_k - a'_k) / \|a_k - a'_k\|_2,$$

где $a'_k = (\pm \|a_k\|_2, 0, \dots, 0)^T$, а знак первой компоненты вектора a'_k для повышения вычислительной устойчивости должен быть противоположен знаку первой компоненты вектора a_k . Напомним, что матрицы Q_k в явном виде не строятся, а определяются векторами нормали w_k .

▷₁ Запишите полностью алгоритм приведения матрицы к форме Хессенберга методом отражений.

Замечание 3.5. Если исходная матрица A симметрична, то нетрудно заметить, что ее форма Хессенберга является симметричной трехдиагональной матрицей. Это обстоятельство позволяет существенно снизить вычислительные затраты.

▷₂ Постройте соответствующий алгоритм.

3.4.3. Преобразование вращения

Для эффективного применения QR -алгоритма к матрице в форме Хессенберга необходимо познакомиться еще с одним видом ортогональных преобразований.

Определение 3.4. Матрицей элементарного вращения называется матрица вида

$$V_{pq}(\theta) = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & c & s & \\ & & & & \ddots & \\ & & & -s & c & \\ & & & & & 1 & \\ & & & & & & \ddots & \\ & & & & & & & 1 \end{bmatrix} \begin{matrix} p \\ q \end{matrix}, \quad (3.11)$$

где $c = \cos \theta$, $s = \sin \theta$. Умножение такой матрицы на вектор $x \in \mathbb{R}^n$ эквивалентно повороту этого вектора на угол θ в плоскости, соответствующей координатам с номерами p и q .

▷₃ Докажите, что матрица вращения ортогональна.

Заметим, что умножение на матрицу V_{pq} изменяет в произвольном $x \in \mathbb{R}^n$ только p -й и q -й элементы:

$$V_{pq} \begin{bmatrix} x_1 \\ \vdots \\ x_p \\ \vdots \\ x_q \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ c x_p + s x_q \\ \vdots \\ -s x_p + c x_q \\ \vdots \\ x_n \end{bmatrix}.$$

Сейчас мы рассмотрим как такие преобразования используются для построения QR -разложения хессенберговой матрицы.

3.4.4. QR -разложение для формы Хессенберга

Пусть A имеет форму Хессенберга. Для приведения ее к верхнетреугольному виду достаточно обнулить $n - 1$ элементов, находящихся под главной диагональю. Причем, как уже говорилось, соответствующие линейные преобразования должны быть ортогональными. Преобразования

вращения как нельзя лучше подойдут для этого:

$$\underbrace{V_{n-1,n} \dots V_{23} V_{12}}_{Q^{-1}} A = R, \quad (3.12)$$

где матрицы $V_{i,i+1} = V_{i,i+1}(\theta_i)$ имеют вид (3.11) и определяются лишь параметрами $s_i = \sin \theta_i$ и $c_i = \cos \theta_i$. Проиллюстрируем это схемой для случая $n = 4$:

$$\begin{aligned} \begin{bmatrix} c_1 & s_1 & & \\ -s_1 & c_1 & & \\ & & 1 & \\ & & & 1 \end{bmatrix} \begin{bmatrix} \times & \times & \times & \times \\ * & \times & \times & \times \\ & * & \times & \times \\ & & * & \times \end{bmatrix} &= \begin{bmatrix} \times & \times & \times & \times \\ & \times & \times & \times \\ & & * & \times & \times \\ & & & * & \times \end{bmatrix}, \\ \begin{bmatrix} 1 & & & \\ & c_2 & s_2 & \\ & -s_2 & c_2 & \\ & & & 1 \end{bmatrix} \begin{bmatrix} \times & \times & \times & \times \\ & \times & \times & \times \\ & & * & \times & \times \\ & & & * & \times \end{bmatrix} &= \begin{bmatrix} \times & \times & \times & \times \\ & \times & \times & \times \\ & & \times & \times \\ & & & * & \times \end{bmatrix}, \\ \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & c_3 & s_3 \\ & & -s_3 & c_3 \end{bmatrix} \begin{bmatrix} \times & \times & \times & \times \\ & \times & \times & \times \\ & & \times & \times \\ & & & * & \times \end{bmatrix} &= \begin{bmatrix} \times & \times & \times & \times \\ & \times & \times & \times \\ & & \times & \times \\ & & & \times & \times \end{bmatrix}. \end{aligned}$$

Таким образом, QR -разложение матрицы Хессенберга определяется матрицей R и набором $\{c_i, s_i\}$, $i = \overline{1, n-1}$.

Согласно (3.12), имеем

$$Q = (V_{n-1,n} \dots V_{23} V_{12})^{-1} = V_{12}^T V_{23}^T \dots V_{n-1,n}^T,$$

поэтому для завершения итерации QR -алгоритма, т. е. для вычисления $A' = RQ$ нужно последовательно умножить полученную матрицу R справа на $V_{i,i+1}^T$ для i от 1 до $n-1$. Каждое такое преобразование будет изменять в матрице R столбцы с номерами i и $i+1$. Важно, что A' также будет иметь форму Хессенберга (иначе изначальное приведение A к данной форме было бы бессмысленным).

▷₄ Докажите это утверждение.

Чтобы полностью определить алгоритм, осталось вывести формулы для параметров $\{c_i, s_i\}$, обеспечивающие обнуление «поддиагональных» элементов в ходе (3.12). Для этого достаточно рассмотреть случай 2×2 :

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \sqrt{\alpha^2 + \beta^2} \\ 0 \end{bmatrix}.$$

Первая компонента правой части равна длине исходного вектора $[\alpha, \beta]^T$ в силу ортогональности преобразования вращения. Второе уравнение этой СЛАУ имеет вид

$$-s\alpha + c\beta = 0,$$

откуда сразу получаем $t = s/c = \beta/\alpha$, или

$$c = \frac{\alpha}{\sqrt{\alpha^2 + \beta^2}}, \quad s = \frac{\beta}{\sqrt{\alpha^2 + \beta^2}}. \quad (3.13)$$

▷₅ Запишите псевдокод одной итерации QR -алгоритма для матрицы в форме Хессенберга. Оцените его сложность.

▷₆ Выполните предыдущее задание для случая $A = A^T$ (трехдиагональная форма Хессенберга).

3.4.5. Детали реализации алгоритма

Итак, с учетом всего вышеизложенного мы можем записать следующее.

Примитивный QR -алгоритм

- 1) Матрица A приводится к форме Хессенберга A_0 .
- 2) Выполняются итерации QR -алгоритма с начальной матрицей A_0 до тех пор, пока не достигается нужной точности.
- 3) Из полученной матрицы $A_N \approx \tilde{R}$ «извлекаются» собственные значения.

Обсудим данный алгоритм. Первый пункт нами рассмотрен (пп. 3.4.2). То, как выполнять QR -итерации с использованием преобразований вращения (пункт 2), мы обсудили в пп. 3.4.4.

Для полной ясности приведем пример «извлечения» собственных значений в пункте 3 (см. основную теорему 3.1). Пусть $n = 4$ и после выполнения необходимого числа итераций QR -алгоритма получена матрица

$$A_N = \begin{bmatrix} \alpha_1 & \times & \times & \times \\ 0 & \alpha_2 & \beta & \times \\ 0 & \gamma & \alpha_3 & \times \\ 0 & 0 & 0 & \alpha_4 \end{bmatrix}, \quad \gamma \neq 0.$$

Видно, что у матрицы два вещественных собственных значения: α_1 и α_4 . Оставшаяся комплексно-сопряженная пара собственных значений

находится как решение характеристического уравнения

$$\begin{vmatrix} \alpha_2 - \lambda & \beta \\ \gamma & \alpha_3 - \lambda \end{vmatrix} = 0.$$

Замечание 3.6. В заключение необходимо отметить, что открытым остался вопрос о критерии остановке итераций в пункте 2. К сожалению, эффективное решение данного вопроса требует внесения в алгоритм дополнительных деталей, обоснование которых выходит за рамки нашего курса. Необходимо понимать также, что приведенная схема алгоритма является весьма упрощенной и недостаточно эффективной. Рациональная реализация QR -алгоритма является намного более сложной.

3.5. Метод вращений Якоби

3.5.1. Симметричная проблема собственных значений

Пусть матрица A — вещественная и симметричная: $A^T = A$. Проблема собственных значений для такой матрицы является более простой, чем в общем случае, поскольку:

- A диагонализируема: существует матрица X такая, что

$$X^{-1}AX = D = \text{diag}(\lambda_1, \dots, \lambda_n), \quad (3.14)$$

причем X — ортогональна;

- все λ_i вещественны.

Также из (3.14) имеем $AX = XD$, т. е. столбцы матрицы X являются собственными векторами A :

$$Ax_i = \lambda_i x_i.$$

3.5.2. Общая схема вращений Якоби

Метод вращений Якоби позволяет вычислить все собственные значения и векторы вещественной симметричной матрицы, т. е. решает полную проблему собственных значений для такой матрицы. Суть метода состоит в построении последовательности матриц

$$A^{(0)} = A, \quad A^{(m)} = X_m^T A^{(m-1)} X_m, \quad m = 1, 2, \dots,$$

где X_m — матрицы элементарного вращения вида (3.11), которые строятся таким образом, что

$$A^{(m)} \xrightarrow{m \rightarrow \infty} D, \quad (3.15)$$

где D — некоторая диагональная матрица. Таким образом, для достаточно большого M получим

$$A^{(M)} \approx \text{diag}(\lambda_1, \dots, \lambda_n), \quad X \approx X^{(M)} = X_1 X_2 \dots X_M.$$

В результате на диагонали матрицы $A^{(M)}$ будут находиться приближенные собственные значения матрицы A , а приближенные собственные векторы будут столбцами матрицы $X^{(M)}$.

Основной вопрос теперь состоит в том, каким образом следует выбирать матрицы вращения X_m , чтобы достичь сходимости к диагональной матрице. Для ответа нам понадобится следующее понятие.

Определение 3.5. *Нормой Фробениуса* называется матричная норма $\|\cdot\|_F$, определяемая как

$$\|A\|_F = \sqrt{\sum_{i,j}^n |a_{ij}|^2}.$$

▷₁ Докажите, что ортогональные преобразования подобия сохраняют норму Фробениуса.

Рассмотрим величины

$$\text{off}(A) = \sum_{i \neq j}^n a_{ij}^2, \quad \text{on}(A) = \sum_{i=1}^n a_{ii}^2.$$

По определению имеем

$$\text{on}(A) + \text{off}(A) = \|A\|_F^2.$$

Условие сходимости (3.15) теперь можно записать в эквивалентном виде

$$\text{off}(A^{(m)}) \xrightarrow{m \rightarrow \infty} 0.$$

Как уже говорилось, матрица X_m , определяющая переход от $A^{(m-1)}$ к $A^{(m)}$, является матрицей элементарного вращения (3.11). Таким образом, $X_m = V_{pq}(\theta)$ определяется тремя параметрами: p , q и θ . Чтобы понять, каким образом следует выбирать эти параметры, рассмотрим, как преобразование подобия

$$A \mapsto A' = V_{pq}^T A V_{pq} \quad (3.16)$$

изменяет матрицу A (см. рисунок).

Видно, что преобразование (3.16) изменяет в матрице A только строки и столбцы с индексами p и q . В частности, диагонали матриц A и A' отличаются лишь элементами на позициях (p, p) и (q, q) . Введем следующие обозначения: $a_{pp} = a$, $a_{qq} = b$, $a'_{pp} = \alpha$, $a'_{qq} = \beta$. Тогда

$$\text{on}(A') = \text{on}(A) - a^2 - b^2 + \alpha^2 + \beta^2. \quad (3.17)$$

$$\begin{array}{c}
\begin{array}{|c|c|c|c|c|}
\hline
1 & & & & \\
\hline
& 1 & & & \\
\hline
& & c & & -s \\
\hline
& & & 1 & \\
\hline
& & & & 1 \\
\hline
& & & & \\
\hline
s & & & & c \\
\hline
& & & & 1 \\
\hline
\end{array}
\end{array}
\begin{array}{c}
p \\
q
\end{array}
\begin{array}{c}
V_{pq}^T \\
A \\
V_{pq}
\end{array}
=
\begin{array}{|c|c|c|c|c|}
\hline
1 & & & & \\
\hline
& 1 & & & \\
\hline
& & c & & s \\
\hline
& & & 1 & \\
\hline
& & & & 1 \\
\hline
& & & & \\
\hline
-s & & & & c \\
\hline
& & & & 1 \\
\hline
\end{array}
\begin{array}{c}
A'
\end{array}$$

Обозначим также $a_{pq} = a_{qp} = e$, $a'_{pq} = a'_{qp} = \varepsilon$. Из (3.16) имеем соотношение

$$\begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{bmatrix} a & e \\ e & b \end{bmatrix} \begin{bmatrix} c & s \\ -s & c \end{bmatrix} = \begin{bmatrix} \alpha & \varepsilon \\ \varepsilon & \beta \end{bmatrix}. \quad (3.18)$$

Поскольку ортогональное преобразование подобия сохраняет норму Фробениуса, получим

$$\alpha^2 + \beta^2 + 2\varepsilon^2 = a^2 + b^2 + 2e^2.$$

Используем это тождество в (3.17):

$$\text{on}(A') = \text{on}(A) + 2(e^2 - \varepsilon^2),$$

что с учетом

$$\|A\|_F = \|A'\|_F$$

равносильно

$$\text{off}(A') = \text{off}(A) - 2(e^2 - \varepsilon^2). \quad (3.19)$$

Из равенства (3.19) видно, что наименьшее значение $\text{off}(A')$ достигается, если величины s и c (т. е. угол θ) выбраны таким образом, что $\varepsilon = a'_{pq} = 0$. Для того чтобы величина $\text{off}(A)$ уменьшилась как можно больше, нужно выбрать в качестве обнуляемого элемента a_{pq} максимальный по модулю недиагональный элемент матрицы A .

Исходя из вышесказанного, получим следующий общий алгоритм.

Базовый алгоритм метода Якоби. Пока $\text{off}(A)$ недостаточно мало:

- 1) среди a_{ij} для $i < j$ выбрать максимальный по модулю a_{pq} ;
- 2) выбрать c и s таким образом, чтобы после преобразования (3.16) получилось $a'_{pq} = 0$;
- 3) $A \leftarrow V_{pq}^T A V_{pq}$.

Согласно (3.19), для этого алгоритма имеем

$$\text{off}(A') = \text{off}(A) - 2e^2. \quad (3.20)$$

Здесь $A' = A^{(m+1)}$, $A = A^{(m)}$, $e = a_{pq}^{(m)}$. Важно понимать, что каждый шаг алгоритма в общем случае «портит» нули, сделанные на предыдущем шаге.

Оценим скорость сходимости алгоритма. Если $a_{pq} = e$ — максимальный по модулю недиагональный элемент матрицы A , то справедлива оценка

$$\text{off}(A) \leq 2Ne^2,$$

где $N = n(n-1)/2$. Тогда из (3.20) получим

$$\text{off}(A') \leq \left(1 - \frac{1}{N}\right) \text{off}(A),$$

что означает

$$\text{off}(A^{(m)}) \leq \left(1 - \frac{1}{N}\right)^m \text{off}(A). \quad (3.21)$$

В частности, при $n = 2$ очевидно имеем сходимость за одну итерацию.

Согласно (3.21), метод вращений Якоби сходится со скоростью геометрической прогрессии. Однако эта оценка слишком груба, и на практике метод сходится быстрее.

3.5.3. Расчетные формулы метода

Для окончательного определения метода осталось вывести формулы, по которым вычисляются элементы матрицы V_{pq} во втором пункте алгоритма. Из (3.18) получим

$$\alpha = a c^2 + b s^2 - 2e sc, \quad (3.22a)$$

$$\beta = a s^2 + b c^2 + 2e sc, \quad (3.22б)$$

$$\varepsilon = e(c^2 - s^2) + (a - b)sc = 0. \quad (3.22в)$$

Для решения уравнения (3.22в) введем переменные

$$t = \text{tg } \theta = \frac{s}{c}, \quad z = \frac{b - a}{2e}.$$

После простых преобразований из (3.22в) получим

$$t^2 + 2zt - 1 = 0,$$

откуда

$$t = -z \pm \sqrt{z^2 + 1}. \quad (3.23)$$

Использование этой формулы на практике приводит к большим ошибкам округления, поэтому ее нужно переписать в более подходящем для машинных вычислений виде. Домножив (3.23) на $z \pm \sqrt{z^2 + 1}$, получим

$$t = \frac{1}{z \pm \sqrt{z^2 + 1}}.$$

Важно выбрать из этих двух корней наименьший по модулю. При $z \neq 0$ он равен

$$t = \frac{\text{sign } z}{|z| + \sqrt{z^2 + 1}}. \quad (3.24)$$

После этого вычислим

$$c = \frac{1}{\sqrt{1 + t^2}}, \quad s = tc. \quad (3.25)$$

Замечание 3.7. В случае $z = 0$ необходимо использовать формулу (3.23), которая дает $t = \pm 1$, или $s = c = 1/\sqrt{2}$. Обратите внимание, что формула (3.24) в этом случае даст неверный результат.

Теоретически теперь нам остается с помощью найденных значений s и c вычислить A' по формуле (3.16). Однако для того чтобы метод был эффективен, вычисления организуются следующим образом.

Как видно на приведенном выше рисунке, преобразование $A \mapsto A'$ заключается в изменении только строк и столбцов с индексами p и q в матрице A . Для диагональных элементов $a'_{pp} = \alpha$ и $a'_{qq} = \beta$ имеют место формулы (3.22а), (3.22б), $a'_{pq} = a'_{qp} = 0$ по построению, а для $j \neq p$, $j \neq q$ из (3.16) получим

$$\begin{aligned} a'_{jp} &= a'_{pj} = ca_{pj} - sa_{qj}, \\ a'_{jq} &= a'_{qj} = sa_{pj} + ca_{qj}. \end{aligned} \quad (3.26)$$

Для вычислительной устойчивости нужно представить вышеперечисленные формулы в виде

$$a'_{ij} = a_{ij} + \delta_{ij},$$

где δ_{ij} — некоторая поправка. Так, из (3.22а), (3.22б) с учетом (3.22в) получим

$$\begin{aligned} a'_{pp} &= a_{pp} - ta_{pq}, \\ a'_{qq} &= a_{qq} + ta_{pq}, \end{aligned} \quad (3.27)$$

а вместо (3.26) имеем

$$\begin{aligned} a'_{pj} &= a_{pj} - s(a_{qj} + \tau a_{pj}), \\ a'_{qj} &= a_{qj} + s(a_{pj} - \tau a_{qj}), \text{ где } \tau = \frac{s}{1+c}. \end{aligned} \quad (3.28)$$

▷₂ Запишите часть вычислительного алгоритма, отвечающую за построение системы собственных векторов матрицы A .

Замечания по практической реализации

- 1) В памяти следует хранить только верхний треугольник матрицы A .
- 2) Матрицы вращения V_{pq} в память записывать не следует, все преобразования осуществляются с использованием только величин c и s .
- 3) Преобразования (3.16) осуществляются по формулам (3.27), (3.28), при этом нужно грамотно учитывать симметрию матрицы.
- 4) При больших n выбор максимального по модулю элемента требует слишком много времени, поэтому, как правило, элементы a_{pq} для обнуления выбирают циклически: $a_{12}, \dots, a_{1n}, a_{23}, \dots, a_{2n}, \dots, a_{n-1,n}$. Обычно для получения решения в пределах машинной точности достаточно 5–6 таких проходов. При этом на первых 2–3 проходах, если модуль a_{pq} достаточно мал, его пропускают.

Глава 4

Методы решения нелинейных уравнений и систем

- 4.1. Решение нелинейных уравнений
- 4.2. Решение систем нелинейных уравнений
- 4.3. Анализ сходимости итерационных процессов
- 4.4. Численные методы оптимизации

4.1. Решение нелинейных уравнений

4.1.1. Введение

Рассмотрим задачу нахождения корней нелинейных уравнений вида

$$f(x) = 0, \quad f : \mathbb{R} \rightarrow \mathbb{R}. \quad (4.1)$$

Определение 4.1. Если $f^{(i)}(x^*) = 0 \forall i = \overline{0, m-1}$, и $f^{(m)}(x^*) \neq 0$, то говорят, что x^* — корень уравнения (4.1) кратности m .

В общем случае задача (4.1) весьма сложна, поэтому перед тем как приступить к ее решению, необходимо провести следующий анализ:

- 1) определить количество корней и, желательно, их кратность;
- 2) *отделить (локализовать) корни* — определить непересекающиеся интервалы, в каждом из которых находится по одному корню.

Чем точнее отделены корни, тем больше шансов на успешное решение задачи. Основное средство, которое мы будем применять для отделения корней — известный из анализа факт: *если непрерывная на отрезке $[a, b]$ функция f принимает на концах этого отрезка значения противоположных знаков, то на $[a, b]$ находится как минимум один корень уравнения $f(x) = 0$.*

4.1.2. Основные теоретические сведения

Все рассматриваемые нами далее методы являются *итерационными*. Это означает, что они заключаются в построении *по единому правилу* последовательности чисел (x_k) таким образом, чтобы $x_k \rightarrow x^*$ при $k \rightarrow \infty$, где x^* — искомый корень уравнения $f(x) = 0$. Большинство этих методов можно будет представить в виде

$$x_{k+1} = \varphi(x_k), \quad k = 0, 1, 2, \dots, \quad (4.2)$$

где x_0 — начальное приближение; $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ — некоторая функция, однозначно определяющая метод. Функция φ , естественно, должна зависеть от f . Переходя к пределу в (4.2), видим, что искомое решение x^* является также и корнем уравнения

$$x = \varphi(x). \quad (4.3)$$

Определение 4.2. Величину $\varepsilon_k = x_k - x^*$ будем называть погрешностью k -й итерации.

Основная характеристика качества метода — скорость, с которой ε_k стремится к 0.

Определение 4.3. Если при всех достаточно больших k имеет место неравенство

$$|\varepsilon_{k+1}| \leq C|\varepsilon_k^p|,$$

где $C < \infty$, $p \in \mathbb{R}$, то говорят, что метод, по которому построена последовательность (x_k) , имеет *порядок (скорость) сходимости p* . Если $p = 1$, то сходимость называют *линейной*, при $1 < p < 2$ — *сверхлинейной*, $p = 2$ — *квадратичной* и т. д.

Для методов вида (4.2) существует относительно простой способ определения порядка сходимости.

Лемма 4.1. Пусть функция φ имеет p непрерывных производных в окрестности корня x^* и

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

Тогда итерационный процесс (4.2) имеет порядок сходимости p и справедлива оценка

$$|\varepsilon_{k+1}| \leq \frac{M_p}{p!} |\varepsilon_k^p|, \quad (4.4)$$

где M_p — максимум $|\varphi^{(p)}(x)|$ в некоторой фиксированной окрестности x^* .

Доказательство. Воспользуемся формулой Тейлора:

$$\begin{aligned} \varepsilon_{k+1} = x_{k+1} - x^* &= \varphi(x_k) - x^* = \varphi(x^* + \varepsilon_k) - x^* = \varphi(x^*) + \\ &+ \varphi'(x^*)\varepsilon_k + \dots + \frac{\varphi^{(p-1)}(x^*)}{(p-1)!} \varepsilon_k^{p-1} + \frac{\varphi^{(p)}(\xi)}{p!} \varepsilon_k^p - x^* = \frac{\varphi^{(p)}(\xi)}{p!} \varepsilon_k^p, \end{aligned}$$

где $\xi \in (x^*, x^* + \varepsilon_k)$. Отсюда получим оценку (4.4). \square

4.1.3. Метод бисекции

Если корень локализован на отрезке $[a, b]$ и $f(a)f(b) < 0$, то наиболее простым способом приближенного вычисления такого корня является метод *бисекции (дихотомии, половинного деления)*. При этом данный метод отличается высокой надежностью: в случае непрерывной функции f он всегда сходится к одному из корней уравнения, расположенных на $[a, b]$.

Алгоритм метода приведен ниже и проиллюстрирован на рис. 4.1.

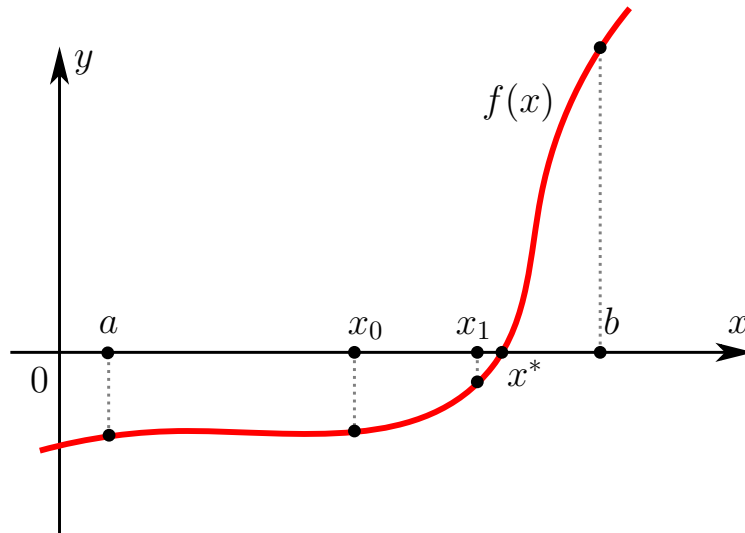


Рис. 4.1. Метод бисекции

Метод бисекции, общая схема

```

1: while  $|b - a| > 2\varepsilon$  do
2:    $x \leftarrow (a + b)/2$ 
3:   if  $f(x)f(a) < 0$  then
4:      $b \leftarrow x$ 
5:   else
6:      $a \leftarrow x$ 
7:   end if
8: end while

```

Здесь ε — точность, с которой требуется найти корень.

▷₁ Определите порядок сходимости метода бисекции.

▷₂ Найдите априорную оценку погрешности $|\varepsilon_k| = |x_k - x^*|$ для метода бисекции.

4.1.4. Метод Ньютона

В случаях, когда известен вид $f'(x)$, метод Ньютона является одним из наиболее эффективных. Рассмотрим некоторое приближение x_k . Приближим график функции f касательной в точке x_k :

$$f(x) \approx y(x) = f(x_k) + (x - x_k)f'(x_k).$$

В качестве следующего приближения к решению уравнения $f(x) = 0$ возьмем корень уравнения $y(x) = 0$ (рис. 4.2):

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, \dots \quad (4.5)$$

Формула (4.5) и определяет метод Ньютона.

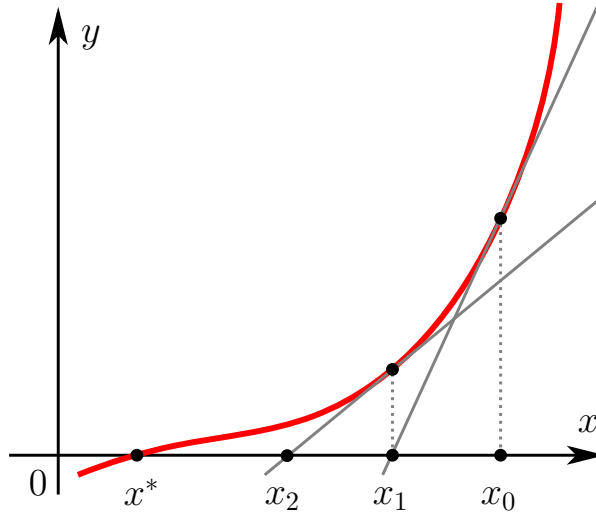


Рис. 4.2. Метод Ньютона

Исследуем порядок сходимости метода Ньютона. Запишем (4.5) в виде

$$x_{k+1} = \varphi(x_k), \quad \varphi(x) = x - \frac{f(x)}{f'(x)}. \quad (4.6)$$

Из (4.6) по лемме 4.1 легко получить порядок сходимости метода Ньютона:

$$\varphi'(x) = 1 - \frac{f'(x)^2 - f(x)f''(x)}{f'(x)^2} = \frac{f(x)f''(x)}{f'(x)^2}, \quad (4.7)$$

откуда при условии $f'(x^*) \neq 0$ получим

$$\varphi'(x^*) = 0. \quad (4.8)$$

Вычисляя $\varphi''(x)$, увидим, что $\varphi''(x^*) \neq 0$, следовательно, если x^* — корень кратности 1, то метод Ньютона имеет второй порядок сходимости.

Предположим теперь, что x^* — корень кратности m . Тогда для достаточно гладкой функции f в его окрестности для точки $x = x^* + \delta$ имеем

$$f(x) = \underbrace{f(x^*) + f'(x^*)\delta + \dots + \frac{f^{(m-1)}(x^*)}{(m-1)!}\delta^{m-1}}_{=0} + O(\delta^m),$$

т. е.

$$f(x) \approx c(x - x^*)^m,$$

где c — некоторая константа. Подставляя это представление в (4.7), получим

$$\varphi'(x) \approx \frac{m-1}{m},$$

следовательно, по формуле (4.4) имеем, что в случае корня кратности m метод Ньютона сходится со скоростью геометрической прогрессии со знаменателем $\frac{m-1}{m}$, т. е. линейно.

4.1.5. Условия сходимости метода Ньютона

Теорема 4.1. Пусть на отрезке $[a, b]$ функция f дважды непрерывно дифференцируема, $f(a)f(b) < 0$, а $f'(x)$ и $f''(x)$ при всех $x \in [a, b]$ не равны нулю и сохраняют знак. Тогда метод Ньютона сходится к решению уравнения $f(x) = 0$ при всех начальных приближениях x_0 , удовлетворяющих условию

$$f(x_0)f''(x_0) > 0. \quad (4.9)$$

Доказательство. По условию теоремы функция f на $[a, b]$ монотонна и не имеет точек перегиба. Изобразив график такой функции нетрудно интуитивно убедиться в справедливости утверждения теоремы: если выбрать начальное приближение с нужной стороны от корня (см. условие (4.9)), то согласно геометрическому смыслу метода Ньютона получим монотонно сходящуюся к корню последовательность.

Чтобы доказать теорему строго, для определенности зафиксируем знак у f' и f'' : пусть $f'(x) > 0$, $f''(x) > 0$ для всех $x \in [a, b]$. Тогда по условию $f(b) > 0$ и согласно (4.9) необходимо выбрать начальное приближение из условия $f(x_0) > 0$, то есть $x_0 \in (x^*, b]$. Теперь нужно показать то, что подсказывает геометрическая интуиция: последовательность метода

Ньютона $x_{k+1} = x_k - f(x_k)/f'(x_k)$ является монотонно убывающей и сходится к x^* . Начнем с x_1 . Согласно сделанным предположениям имеем $x_1 < x_0$, но важно убедиться, что мы не «перескочили» на противоположную сторону от корня, то есть что $f(x_1) > 0$. Для этого представим $f(x_1)$ по формуле Тейлора:

$$\begin{aligned} f(x_1) &= f\left(x_0 - \frac{f(x_0)}{f'(x_0)}\right) = \\ &= f(x_0) - f'(x_0) \frac{f(x_0)}{f'(x_0)} + \frac{f''(\xi)}{2} \left(\frac{f(x_0)}{f'(x_0)}\right)^2 > 0. \end{aligned}$$

Эти же самые рассуждения останутся в силе, если мы заменим теперь x_0 и x_1 на x_k и x_{k+1} . То есть, при всех k действительно имеем

$$x^* < x_k < x_{k-1} < \dots < x_1 < x_0. \quad (4.10)$$

Монотонно убывающая и ограниченная снизу последовательность $\{x_k\}$ обязана иметь предел (теорема Вейерштрасса), который пока что обозначим \tilde{x} . Чтобы убедиться, что $\tilde{x} = x^*$ достаточно перейти к пределу при $k \rightarrow \infty$ в формуле (4.5), задающей метод Ньютона. \square

▷₃ Проверьте доказательство теоремы в случае, когда f' и f'' имеют знаки, отличные от рассмотренных нами. Попробуйте доказать теорему «в общем», без рассмотрения всех частных случаев.

Практическая оценка погрешности

При практической реализации любых итерационных методов хочется иметь критерии остановки итераций, которые бы гарантировали нахождение решения с требуемой точностью. Такие критерии, как правило, доступны лишь в простейших ситуациях. В общем случае о сходимости процесса можно судить по величине разности двух соседних приближений $x_k - x_{k-1}$, либо по величине невязки $f(x_k)$. Однако необходимо понимать, что зачастую эти величины могут иметь мало общего с интересующей нас погрешностью $x_k - x^*$. Следующий результат [10] позволяет установить связь между этими величинами.

Теорема 4.2. Пусть выполнены условия теоремы 4.1. Тогда погрешность метода Ньютона $\varepsilon_k = x_k - x^*$ удовлетворяет условиям

$$|x_k - x^*| \leq \frac{|f(x_k)|}{m_1}, \quad (4.11)$$

и

$$|x_k - x^*| \leq \frac{M_2}{2m_1}(x_k - x_{k-1})^2, \quad (4.12)$$

где

$$m_1 = \min_{x \in [a, b]} |f'(x)|, \quad M_2 = \max_{x \in [a, b]} |f''(x)|. \quad (4.13)$$

Доказательство. По теореме Лагранжа о среднем значении на отрезке $[x^*, x_k]$ имеем

$$f(x_k) = f(x^*) + (x_k - x^*)f'(\xi), \quad \xi \in (x^*, x_k),$$

откуда сразу получаем (4.11). Для вывода второй оценки применим формулу Тейлора:

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

где $\eta \in (x_{k-1}, x_k)$. Расписывая x_k по определению метода Ньютона, отсюда легко получаем

$$|f(x_k)| \leq \frac{1}{2}M_2(x_k - x_{k-1})^2. \quad (4.14)$$

Теперь осталось лишь подставить это неравенство в (4.11). \square

4.1.6. Модификации метода Ньютона

Метод Ньютона с постоянной производной

В случаях, когда производная $f'(x)$ не известна или ее слишком дорого вычислять, можно в методе Ньютона (4.5) заменить $f'(x_k)$ на какую-то константу μ :

$$x_{k+1} = x_k - \frac{f(x_k)}{\mu}, \quad k = 0, 1, \dots \quad (4.15)$$

Очевидный выбор — $\mu \approx f'(x_0)$. Геометрическая интерпретация метода в данном случае изображена на рис. 4.3.

▷₄ Определите порядок сходимости метода.

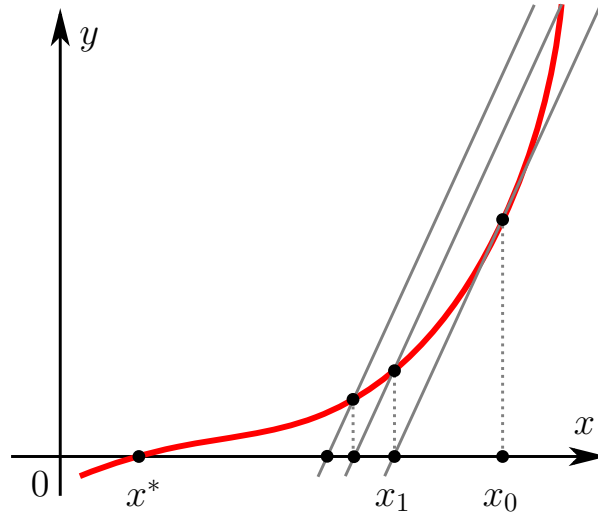


Рис. 4.3. Метод Ньютона с постоянной производной

Метод секущих

Наиболее естественное видоизменение метода Ньютона состоит в замене касательной прямой на секущую при геометрическом построении метода.

Пусть имеется два приближения к корню: x_k и x_{k-1} . Составим уравнение секущей для функции f по этим точкам:

$$y(x) = f(x_k) + (x - x_k) \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}.$$

Положим x_{k+1} равным корню уравнения $y(x) = 0$:

$$x_{k+1} = x_k - f(x_k) \frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})}, \quad k = 0, 1, \dots \quad (4.16)$$

Это — метод секущих. Его геометрическая интерпретация изображена на рис. 4.4. Этот метод нельзя представить в виде (4.2), поэтому порядок сходимости метода вычисляется нетривиально. Он равен золотому сечению:

$$m = \frac{1 + \sqrt{5}}{2} \approx 1,618,$$

т. е. метод секущих сходится сверхлинейно. В качестве начальных приближений x_1 и x_0 можно брать концы отрезка $[a, b]$, на котором локализован корень.

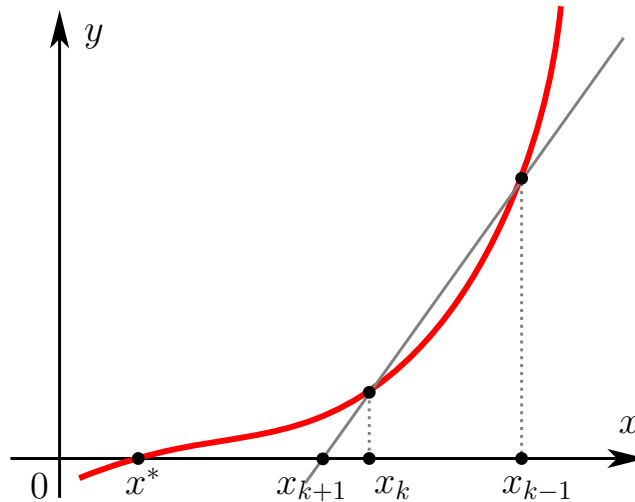


Рис. 4.4. Метод секущих

Метод хорд

У метода секущих есть недостаток: члены последовательности x_k могут выходить за пределы отрезка локализации $[a, b]$, что может привести к расходимости итерационного процесса.

▷₅ Приведите пример, когда при $x_0 = a$ и $x_1 = b$ приближение x_3 , построенное по методу секущих, не лежит в $[a, b]$.

Избежать этой проблемы можно, зафиксировав один из концов секущей линии:

$$x_{k+1} = x_k - f(x_k) \frac{x_k - x_0}{f(x_k) - f(x_0)}, \quad (4.17)$$

причем в качестве x_0 нужно брать тот конец отрезка $[a, b]$, в котором знаки f и f'' совпадают.

▷₆ Обоснуйте это правило.

▷₇ Определите порядок сходимости метода хорд.

Геометрический смысл метода хорд (4.17) изображен на рис. 4.5.

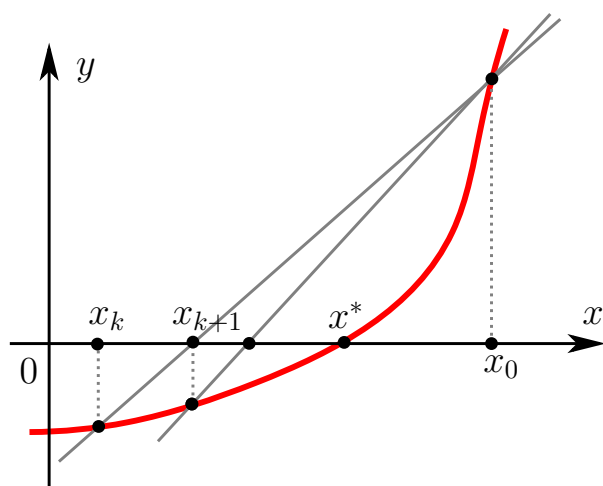


Рис. 4.5. Метод хорд

$$x^{k+1} = x^k - f'(x^k)^{-1} f(x^k), \quad (4.21)$$

где $f' = \frac{\partial f}{\partial x}$ — матрица Якоби,

$$f' = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \cdots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}. \quad (4.22)$$

Для обоснования формулы (4.21) рассмотрим $x \in \mathbb{R}^n$, приращение $\Delta x \in \mathbb{R}^n$ и воспользуемся формулой Тейлора для функции нескольких переменных:

$$f_i(x + \Delta x) = f_i(x) + \sum_{j=1}^n \frac{\partial f_i}{\partial x_j}(x) \Delta x_j + O(\|\Delta x\|^2), \quad i = \overline{1, n}.$$

В векторной форме эти равенства можно записать как

$$f(x + \Delta x) = f(x) + f'(x)\Delta x + O(\|\Delta x\|^2). \quad (4.23)$$

Пусть x^k — некоторое приближение к x^* . В идеале нам нужно найти такую поправку Δx^* , для которой $x^k + \Delta x^* = x^*$, или

$$f(x^k + \Delta x^*) = 0.$$

Отбрасывая остаточный член в формуле (4.23), имеем

$$0 = f(x^k + \Delta x^*) \approx f(x^k) + f'(x^k)\Delta x^* \Rightarrow \Delta x^* \approx -f'(x^k)^{-1}f(x^k) = \Delta x^k.$$

Вычисляя

$$x^* \approx x^{k+1} = x^k + \Delta x^k,$$

получим формулу (4.21).

Замечание 4.1. Важно понимать, что формула (4.21) в явном виде практически никогда не используется для расчетов, так как для этого нужно иметь аналитический вид матрицы $f'(x)^{-1}$.

Вместо этого каждую итерацию метода Ньютона реализуют в два этапа:

1) сначала находят поправку Δx^k как решение СЛАУ

$$f'(x^k)\Delta x^k = -f(x^k); \quad (4.24a)$$

2) затем вычисляют

$$x^{k+1} = x^k + \Delta x^k. \quad (4.24б)$$

▷₁ Оцените сложность реализации метода Ньютона по формулам (4.24).

При больших n метод Ньютона становится слишком трудоемким. Кроме этого, он не применим, если нет возможности вычислить f' , что довольно часто случается на практике. Для таких случаев разработан ряд модификаций метода Ньютона, направленных на снижение вычислительных затрат.

4.2.3. Метод Ньютона с постоянным якобианом

Для того чтобы снизить затраты на решение СЛАУ (4.24a), рассматривают метод Ньютона с «замороженной» матрицей Якоби:

$$x^{k+1} = x^k - J_0^{-1} f(x_k), \quad (4.25)$$

где $J_0 = f'(x_0)$.

Для реализации такого метода достаточно один раз построить LU -разложение матрицы J_0 , а затем с помощью него решать СЛАУ вида

$$J_0 \Delta x^k = -f(x_k)$$

для нахождения поправок за $O(n^2)$ операций. Данная модификация на порядок снижает вычислительные затраты, но при этом, конечно, снижается скорость сходимости метода.

4.2.4. Обобщения метода секущих

Метод секущих (4.16) можно рассматривать с двух сторон.

1) С одной стороны, он может быть получен путем модификации метода Ньютона с помощью замены

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

т. е. метод секущих основан на приближенном вычислении f' .

2) С другой стороны, данный метод определяет уравнение секущей $y(x)$, которая пересекает график f в точках x_k и x_{k-1} (именно так мы этот метод строили в предыдущем пункте). С этой точки зрения мы строим приближение для f .

В одномерном случае два описанных подхода дают один и тот же метод, однако в многомерном случае мы получим два разных семейства методов.

Первый способ: приближенное вычисление якобиана. Рассмотрим набор приращений $\{h_1, \dots, h_n\}$, $h_i \in \mathbb{R}$, и определим приближения для частных производных f по формуле

$$\frac{\partial f_i}{\partial x_j}(x) \approx \frac{f_i(x_1, \dots, x_j + h_j, \dots, x_n) - f_i(x_1, \dots, x_n)}{h_j} = \psi_{ij}(x), \quad (4.26)$$

откуда имеем $f'(x) \approx \Psi(x)$. Соответствующий метод имеет вид

$$x^{k+1} = x^k - \Psi(x^k)^{-1} f(x^k). \quad (4.27)$$

▷₂ Укажите достоинства и недостатки этого метода.

Кроме (4.26) можно рассматривать и другие варианты приближенного вычисления $f'(x)$.

Второй способ приводит к методу Бройдена.

4.2.5. Метод Бройдена

Обобщением понятия прямой, задаваемой уравнением $y = ax + b$, для пространства \mathbb{R}^n является множество пар $(x, y) \in \mathbb{R}^n \times \mathbb{R}^n$, таких, что

$$y = Ax + b, \quad (4.28)$$

где x, y, b — векторы из \mathbb{R}^n ; A — квадратная матрица.

Пусть нам известны два приближения к x^* : x^k и x^{k-1} . Тогда секущая вида (4.28) определяется соотношениями

$$\begin{cases} Ax^{k-1} + b = f(x^{k-1}), \\ Ax^k + b = f(x^k). \end{cases} \quad (4.29)$$

Очевидно, что эти условия не определяют A и b однозначно, но пока предположим, что мы каким-то образом нашли A . Тогда $f(x) \approx Ax + b$ и следующее приближение x^{k+1} строится из соотношения

$$Ax^{k+1} + b = 0 \quad \Rightarrow \quad x^{k+1} = -A^{-1}b,$$

откуда, выражая из (4.29) $b = f(x^k) - Ax^k$, получим знакомую формулу

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

Общая схема итерации многомерного метода секущих

- 1) По двум предыдущим приближениям x^k и x^{k-1} вычисляется матрица $A = A_k$, определяющая уравнение секущей.
- 2) Находится новое приближение по формуле

$$x^{k+1} = x^k - A_k^{-1}f(x^k). \quad (4.30)$$

Рассмотрим теперь главный вопрос: как определять матрицу A из условий (4.29)? Это можно делать по-разному. Один из наиболее удачных способов был предложен Ч. Бройденом в 1965 г.

Обозначим $\Delta x = x^k - x^{k-1}$, $\Delta f = f(x^k) - f(x^{k-1})$. Из (4.29) имеем

$$A_k \Delta x = \Delta f. \quad (4.31)$$

Основная идея заключается в том, чтобы представить A_k в виде

$$A_k = A_{k-1} + u \Delta x^T, \quad (4.32)$$

где $u \in \mathbb{R}^n$ — вектор неизвестных параметров.

Подставляя (4.32) в (4.31), получим

$$(A_{k-1} + u \Delta x^T) \Delta x = \Delta f,$$

откуда

$$u = \frac{1}{\|\Delta x\|_2^2} (\Delta f - A_{k-1} \Delta x). \quad (4.33)$$

Здесь уже можно остановиться: по u находим A_k , после чего найдем x^{k+1} по формуле (4.30). При этом необходимо решить СЛАУ вида

$$A_k(x^{k+1} - x^k) = -f(x^k). \quad (4.34)$$

При определенной схеме организации вычислений решение таких систем можно получить за $O(n^2)$ операций, но описание такого алгоритма не входит в рамки нашего курса [19].

В самом начале процесса для получения приближения x^1 можно использовать метод Ньютона или метод (4.27). Тогда матрицей A_1 будет матрица $f'(x^0)$ или $\Psi(x^0)$ соответственно.

4.2.6. Другие методы

В заключение рассмотрим *нелинейный метод Гаусса – Зейделя*. По аналогии со случаем СЛАУ здесь в цикле по i от 1 до n по очереди вычисляются компоненты x_i^{k+1} путем решения скалярного уравнения вида

$$f_i(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i, x_{i+1}^k, \dots, x_n^k) = 0, \quad (4.35)$$

относительно x_i . После этого полагается $x_i^{k+1} = x_i$. Для решения скалярных уравнений (4.35) могут быть использованы соответствующие методы, рассмотренные ранее.

▷₃ Постройте нелинейные аналоги методов Якоби и релаксации.

4.3. Анализ сходимости итерационных процессов

Как уже не раз упоминалось, многие из рассмотренных нами методов решения нелинейных уравнений и систем представимы в виде $x^{k+1} = \varphi(x^k)$. Рассмотрим классический способ доказательства сходимости таких итерационных процессов.

4.3.1. Принцип сжимающих отображений

Расстояние между двумя точками $x, y \in \mathbb{R}^n$ условимся обозначать

$$\rho(x, y) = \|x - y\|. \quad (4.36)$$

Определение 4.4. Если для функции $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ существует константа $L < \infty$, такая, что

$$\rho(\varphi(x), \varphi(y)) \leq L\rho(x, y), \quad \forall x, y \in D \subset \mathbb{R}^n,$$

то говорят, что φ удовлетворяет *условию Липшица* (или просто *липшицева*) на множестве D . Число L называют *константой Липшица*.

Определение 4.5. *Сжимающим отображением* называют функцию φ , удовлетворяющую условию Липшица на \mathbb{R}^n с константой $L < 1$. Если же это условие выполнено лишь на некоторой области $D \subset \mathbb{R}^n$, то такое отображение будем называть *локально-сжимающим*.

Схема действия сжимающего отображения $\mathbb{R}^2 \rightarrow \mathbb{R}^2$ изображена на рис. 4.6.

Определение 4.6. Последовательность точек (x^k) , $x^k \in \mathbb{R}^n$, называется *фундаментальной* или *последовательностью Коши*, если

$$\rho(x^k, x^m) \rightarrow 0 \quad \text{при} \quad k, m \rightarrow \infty.$$

Согласно критерию Коши, последовательность (x^k) сходится тогда и только тогда, когда она фундаментальна.

Теорема 4.3 (принцип сжимающих отображений). Пусть $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ — сжимающее отображение. Тогда уравнение

$$\varphi(x) = x \quad (4.37)$$

Докажем единственность. Пусть $\exists x^{**} \neq x^*$ такой, что $\varphi(x^{**}) = x^{**}$. Это приводит к противоречию:

$$\rho(x^*, x^{**}) = \rho(\varphi(x^*), \varphi(x^{**})) \leq L\rho(x^*, x^{**}) < \rho(x^*, x^{**}). \quad \square$$

Следствие 4.1. Устремляя $m \rightarrow \infty$ в формуле (4.39), получим априорную оценку погрешности:

$$\rho(x^k, x^*) \leq \frac{L^k}{1-L} \rho(x^0, x^1). \quad (4.40)$$

▷₁ Получите из формулы (4.40) оценку для количества итераций N , достаточного для получения решения с погрешностью, не превышающей ϵ .

Кроме оценки (4.40) можно получить и *апостериорную* оценку погрешности, которая, как правило, более точна, но для получения которой необходимо проделать k итераций. Положив в (4.40) $x^0 = x^{k-1}$, имеем

$$\rho(x^k, x^*) \leq \frac{L}{1-L} \rho(x^{k-1}, x^k). \quad (4.41)$$

Замечание 4.2. Функция ρ (метрика), которая измеряет расстояние между двумя точками, не обязательно должна иметь вид (4.36). Достаточно лишь, чтобы она удовлетворяла аксиомам метрики:

- 1) $\rho(x, y) = 0 \Leftrightarrow x = y$;
- 2) $\rho(x, y) = \rho(y, x)$;
- 3) $\rho(x, y) \leq \rho(x, z) + \rho(z, y)$.

Замечание 4.3. Принцип сжимающих отображений играет важную роль в функциональном анализе и других разделах математики. Он может быть сформулирован не только для \mathbb{R}^n , но и для произвольных метрических пространств, обладающих свойством *полноты*. В частности, с его помощью легко доказывается существование решений для обыкновенных дифференциальных уравнений, интегральных уравнений и т. д.

4.3.2. Локальный принцип сжимающих отображений

Рассмотрим теперь случай, когда отображение φ является лишь локально-сжимающим.

Определение 4.7. Гипершар радиусом r с центром в точке $x \in \mathbb{R}^n$ будем обозначать

$$B(x, r) = \{y \in \mathbb{R}^n \mid \rho(x, y) \leq r\}.$$

Лемма 4.2. Пусть отображение φ имеет неподвижную точку $x^* = \varphi(x^*)$ и является локально-сжимающим в некоторой ее окрестности $B^* = B(x^*, r)$. Тогда для всех $x_0 \in B^*$ итерационный процесс (4.38) сходится к x^* .

▷₂ Докажите лемму.

Для применения леммы 4.2 необходимо обладать достаточно точной информацией о местоположении корня x^* . Следующий признак сходимости опирается лишь на информацию о φ в окрестности начального приближения x^0 .

Теорема 4.4. Пусть отображение φ является локально-сжимающим на множестве $B = B(x, r)$ с константой $L < 1$. Если имеет место неравенство

$$\rho(x, \varphi(x)) \leq (1 - L)r, \quad (4.42)$$

то x^* — искомый корень уравнения $\varphi(x) = x$ — лежит в B и итерационный процесс (4.38) сходится к x^* для всех $x_0 \in B$.

Доказательство. Если мы докажем, что $\varphi(y) \in B$ при любых $y \in B$, то по общему принципу сжимающих отображений из этого автоматически будет следовать утверждение теоремы. Итак, пусть $y \in B$, тогда

$$\begin{aligned} \rho(\varphi(y), x) &\leq \rho(\varphi(y), \varphi(x)) + \rho(x, \varphi(x)) \leq \\ &\leq L\rho(y, x) + \rho(x, \varphi(x)) \leq Lr + (1 - L)r = r. \end{aligned} \quad \square$$

4.3.3. Применение принципа сжимающих отображений

Рассмотрим случай функции одной переменной (одного нелинейного уравнения).

Теорема 4.5 (Лагранжа о среднем значении). Если функция $\varphi : [a, b] \rightarrow \mathbb{R}$ непрерывна на $[a, b]$ и дифференцируема на (a, b) , то $\exists \xi \in (a, b)$ такое, что

$$\varphi(b) - \varphi(a) = \varphi'(\xi)(b - a).$$

В случае $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ существует простое достаточное условие липшицевости: если φ достаточно гладкая, то согласно теореме Лагранжа имеем

$$|\varphi(x) - \varphi(y)| = |\varphi'(\xi)| \cdot |x - y|, \quad \xi \in (x, y),$$

поэтому если

$$|\varphi'(x)| \leq M \quad \forall x \in [a, b],$$

то φ — липшицева на $[a, b]$ с константой $L = M$.

Таким образом, получаем **признак сходимости 0**: если

$$|\varphi'(x)| \leq M < 1 \quad \forall x \in \mathbb{R}, \quad (4.43)$$

то итерационный процесс (4.38) сходится.

Очевидно, что ограничение вида (4.43) является слишком сильным. Как правило, мы можем рассчитывать на выполнение условий такого рода лишь на некотором подмножестве \mathbb{R} . Лемма 4.2 дает **признак сходимости 1**: если

$$|\varphi'(x)| \leq M < 1 \quad \forall x \in B^* = [x^* - r, x^* + r], \quad (4.44)$$

то итерационный процесс (4.38) сходится к x^* при всех $x_0 \in B^*$.

Из теоремы 4.4 получаем также **признак сходимости 2**: если

$$|\varphi'(x)| \leq M < 1 \quad \forall x \in B = [\xi - r, \xi + r]$$

и при этом

$$|\xi - \varphi(\xi)| \leq (1 - M)r,$$

то итерационный процесс (4.38) сходится к $x^* = \varphi(x^*)$ при всех $x_0 \in B$.

Важно понимать, что все рассмотренные выше признаки являются лишь *достаточными* условиями сходимости.

Применение к методу Ньютона

Напомним, что метод Ньютона для решения уравнения $f(x) = 0$ представляет собой метод типа (4.38), где

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

и

$$\varphi'(x) = \frac{f(x)f''(x)}{f'(x)^2}.$$

Пусть $x^* \in [a, b]$, производные $f'(x)$ и $f''(x)$ непрерывны и не обращаются в нуль при всех $x \in [a, b]$. Тогда $\varphi'(x)$ непрерывна и из тождества

$\varphi'(x^*) = 0$ следует, что существует целая окрестность $B^* = B(x^*, r)$ такая, что

$$|\varphi'(x)| \leq M < 1 \quad \forall x \in B^*.$$

Следовательно, по признаку сходимости 1 метод Ньютона будет сходиться при всех x_0 из B^* .

4.4. Численные методы оптимизации

4.4.1. Введение

Рассмотрим функционал (скалярную функцию векторного аргумента)

$$f : \mathbb{R}^n \rightarrow \mathbb{R}.$$

Задача минимизации функционала является одной из важнейших в прикладной математике: требуется найти такой вектор $x^* \in \mathbb{R}^n$, на котором f принимает наименьшее значение. Это, конечно, не совсем четкая постановка задачи, так как понятие «наименьшее значение» можно понимать по-разному. Можно требовать, чтобы $f(x^*)$ было наименьшим лишь в некоторой открытой окрестности x^* :

$$f(x^*) \leq f(x) \quad \forall x \in U(x^*). \quad (4.45)$$

В этом случае x^* является *точкой локального минимума*, или *локальным минимизатором*³ функционала f . Именно такую задачу мы будем рассматривать в дальнейшем. Гораздо более сложно найти *глобальный* минимизатор

$$x^* = \operatorname{argmin}_{x \in \mathbb{R}^n} f(x), \quad (4.46)$$

так как в общем случае f может иметь сколько угодно локальных минимумов. Такие задачи как правило решаются уже вероятностными методами.

В соответствии с устоявшейся терминологией в дальнейшем мы иногда будем называть f *целевой функцией*. Будем также предполагать, что эта функция дважды непрерывно дифференцируема (частные производные $\partial_{ij} f$ существуют и непрерывны).

4.4.2. Условия локального минимума

Приведенные ниже утверждения рассматривались в курсе математического анализа, поэтому мы изложим их в сжатом и неформальном виде.

³Термин «минимизатор» можно считать калькой с английского, так как вместо него в русскоязычной литературе традиционно используется выражение «точка минимума». Мы будем использовать оба варианта. Во избежание возможных недоразумений проговорим отдельно: минимум функционала f это число $f(x^*) \in \mathbb{R}$, а минимизатор (точка минимума) — сам вектор $x^* \in \mathbb{R}^n$.

Вспомним скалярный случай: если $f : \mathbb{R} \rightarrow \mathbb{R}$ достаточно гладкая, то для того, чтобы $x^* \in \mathbb{R}$ являлась локальным минимизатором, необходимо выполнение условия $f'(x^*) = 0$. Чтобы отличить локальный минимум от максимума и перегиба нужно проверить условие $f''(x^*) > 0$ (при переходе через точку x^* первая производная должна возрастать). Таким образом, условия $f'(x^*) = 0$ и $f''(x^*) > 0$ являются достаточными для того, чтобы объявить x^* локальным минимизатором.

Что же произойдет в случае многих переменных? По-прежнему необходимым условием остается обращение в нуль «производной» функционала f — градиента:

$$\nabla f(x^*) = \left(\partial_1 f(x^*), \dots, \partial_n f(x^*) \right)^T = 0. \quad (4.47)$$

Чтобы сделать это условие достаточным, нужно по аналогии с одномерным случаем потребовать, чтобы при движении точки x из x^* в любом направлении значения $f(x)$ возрастали. Математически это требование формулируется в виде

$$\nabla^2 f(x^*) > 0, \quad (4.48)$$

где $\nabla^2 f$ — «вторая производная» целевой функции, матрица частных производных, или матрица Гессе:

$$\nabla^2 f = (\nabla f)' = \begin{bmatrix} \partial_{11}f & \partial_{12}f & \cdots & \partial_{1n}f \\ \partial_{21}f & \partial_{22}f & \cdots & \partial_{2n}f \\ \vdots & \vdots & \ddots & \vdots \\ \partial_{n1}f & \partial_{n2}f & \cdots & \partial_{nn}f \end{bmatrix}. \quad (4.49)$$

В силу гладкости f эта матрица симметрична (порядок взятия частных производных не имеет значения). С условием положительной определенности (4.48) мы уже встречались. Оно означает, что для любого ненулевого вектора $p \in \mathbb{R}^n$ справедливо неравенство

$$p^T \nabla^2 f(x^*) p > 0. \quad (4.50)$$

Давайте попробуем разобраться откуда взялось это условие. Рассмотрим выражение $f(x + p)$. Если аккуратно разложить его в ряды Тейлора

по каждой переменной и сгруппировать слагаемые, то получится формула, на которой держатся все наши последующие построения — формула Тейлора для функции нескольких переменных:

$$f(x + p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x) p + O(\|p\|^3). \quad (4.51)$$

Если положить $x = x^*$, то второе слагаемое в формуле обнуляется и становится понятно, что при всех p с достаточно маленькой нормой знак выражения $f(x^* + p) - f(x^*)$ совпадает со знаком $p^T \nabla^2 f(x^*) p$. Следовательно, добавление условия (4.48) к условию $\nabla f(x^*) = 0$ становится достаточным для того, чтобы x^* был локальным минимизатором. Проговорим еще раз этот результат.

Теорема 4.6 (Достаточное условие локального минимума). *Если целевая функция f дважды непрерывно дифференцируема, $\nabla f(x^*) = 0$ и $\nabla^2 f(x^*) > 0$, то вектор x^* является точкой локального минимума f .*

4.4.3. Общий вид методов спуска

- 1) Выбрать начальное приближение x^0 . Положить $x \leftarrow x^0$.
- 2) Выбрать направление спуска p с единичной нормой: $\|p\|_2 = 1$.
- 3) Выбрать длину шага α .
- 4) Выполнить спуск: $x \leftarrow x + \alpha p$
- 5) Если не выполнено условие остановки (не достигнута нужная точность), перейти к шагу 2.

Наиболее важными этапами в описанной схеме являются пункты 2 и 3. Начнем с выбора подходящего направления спуска в п. 2.

4.4.4. Выбор направления спуска

Классическое направление спуска — *антиградиент*:

$$p = -\nabla f(x) / \|\nabla f(x)\|_2. \quad (4.52)$$

Получаемый в результате метод называется *методом наискорейшего спуска*⁴. Давайте разберемся почему берется именно такое направление.

⁴Иногда этот метод называют методом градиентного спуска, но правильнее было бы говорить, конечно, «метод антиградиентного спуска».

Рассмотрим поведение целевой функции f в окрестности некоторой точки $x \in \mathbb{R}^n$. Начнем мысленно прямолинейное движение из этой точки по направлению, заданному каким-нибудь вектором $p \in \mathbb{R}^n$ с единичной евклидовой нормой. Как при этом будут изменяться значения функции f ? Как только мы зафиксировали направление p , вместо функции f от n переменных возникает скалярная функция от параметра $\alpha \geq 0$:

$$\Phi(\alpha) = f(x + \alpha p).$$

При решении задачи минимизации мы заинтересованы в выборе направления p , которое бы приводило к наискорейшему убыванию f . Рассмотрим *производную по направлению* p в точке x , которая определяется как

$$\partial_p f(x) = \left. \frac{d}{d\alpha} f(x + \alpha p) \right|_{\alpha=0} = \Phi'(0).$$

Так как

$$\begin{aligned} \Phi'(\alpha) &= \frac{d}{d\alpha} f(x + \alpha p) = \frac{d}{d\alpha} f \left(\begin{bmatrix} x_1 + \alpha p_1 \\ \vdots \\ x_n + \alpha p_n \end{bmatrix} \right) = \\ &= \partial_1 f(x + \alpha p) p_1 + \dots + \partial_n f(x + \alpha p) p_n = \nabla f(x + \alpha p)^T p, \end{aligned} \quad (4.53)$$

то

$$\partial_p f(x) = \nabla f(x)^T p = \|\nabla f(x)\| \cos \theta. \quad (4.54)$$

Здесь θ — угол между векторами $\nabla f(x)$ и p (напомним, что по условию $\|p\|_2 = 1$).

Теперь мы готовы обосновать выбор (4.52): нам нужно направление, в котором f убывает быстрее всего, то есть при котором $\Phi'(0) = \partial_p f(x)$ принимает наименьшее значение. Согласно (4.54), это наименьшее значение достигается при $\cos \theta = -1$, то есть при p , направленном в противоположном от градиента направлении. Заодно мы видим главное свойство градиента: он указывает направление наискорейшего возрастания функции f .

Замечание 4.4. Направление не обязательно должно быть антиградиентом. Достаточно, чтобы оно было направлением спуска, то есть направлением убывания f , которое достигается при $\pi/2 < \theta < 3\pi/2$.

4.4.5. Выбор шага спуска: условия Вольфе

Выберем константу $c_1 \in (0, 1)$ и рассмотрим прямую, задаваемую уравнением

$$l(\alpha) = \phi(0) + c_1 \phi'(0)\alpha = f(x) + (c_1 \nabla f(x)^T p)\alpha. \quad (4.55)$$

Условие достаточного убывания:

$$\phi(\alpha) \leq l(\alpha). \quad (4.56)$$

Условие кривизны:

$$\phi'(\alpha) > c_2 \phi'(0), \quad (4.57)$$

где c_2 — еще один параметр, который выбирается из отрезка $(c_1, 1)$. Учитывая (4.53) это условие можно записать в эквивалентном, но более громоздком виде

$$\nabla f(x + \alpha p)^T p > c_2 \nabla f(x)^T p. \quad (4.58)$$

Условия достаточного убывания (4.56) и кривизны (4.57) вместе называются *условиями Вольфе*.

Заметим, что условие кривизны (4.57) допускает сколь угодно большие положительные значения $\phi'(\alpha)$, что не очень логично, так как в таком случае α может находиться далеко от точки локального минимума функции ϕ . Поэтому вместо него иногда используется *усиленное условие кривизны*:

$$|\phi'(\alpha)| \leq c_2 |\phi'(0)|. \quad (4.59)$$

Это условие не дает функции ϕ в точке α слишком сильно возрастать и гарантирует, что α находится вблизи точки локального минимума ϕ .

При практической реализации простейшая стратегия выбора шага заключается в следующем. Сначала выбирается достаточно большое значение $\alpha = \alpha_{\max}$. Если для него не выполняются условия Вольфе, шаг домножается на некоторую константу $\delta < 1$ и вновь выполняется проверка. Процедура выполняется до тех пор, пока условия Вольфе не будут удовлетворены.

4.4.6. Метод Ньютона

Метод наискорейшего спуска имеет первый порядок сходимости, то есть погрешность $\epsilon^k = x^k - x^*$ уменьшается по линейному закону:

$$\|\epsilon^{k+1}\| \leq C \|\epsilon^k\| \quad (4.60)$$

для некоторого $C \in (0, 1)$ и всех достаточно больших k .

Чтобы достичь более высокой скорости сходимости придется заплатить вычислительными затратами. *Метод Ньютона* для поиска локального минимума представляет собой, по сути, процесс решения системы нелинейных уравнений

$$\nabla f(x) = 0 \quad (4.61)$$

методом Ньютона, который мы изучали ранее:

$$x \leftarrow x - (\nabla^2 f(x))^{-1} \nabla f(x). \quad (4.62)$$

Обратим внимание, что этот метод также можно рассматривать как метод спуска, в котором направление задается формулой

$$p = -(\nabla^2 f(x))^{-1} \nabla f(x), \quad (4.63)$$

а шаг α равен 1. В некоторых случаях значение α в методе Ньютона может быть и отлично от 1.

При выполнении определенных условий метод Ньютона имеет квадратичную сходимость. Основная трудность при его реализации — необходимость вычисления и обращения⁵ матрицы Гессе $\nabla^2 f(x)$. Также существует вероятность сходимости к локальному максимуму вместо минимума, что необходимо учитывать при программной реализации данного метода.

⁵На самом деле, конечно, матрицу не обращают, а решают с ней СЛАУ.

Глава 5

Приближение функций

- 5.1. Интерполяция функций. Интерполяционный многочлен
- 5.2. Интерполяционный многочлен в форме Ньютона
- 5.3. Остаток интерполяции
- 5.4. Сплаины
- 5.5. Приближение кривых
- 5.6. Среднеквадратичные приближения
- 5.7. Приближение поверхностей

5.1. Интерполяция функций.

Интерполяционный многочлен

5.1.1. Общая постановка задачи интерполяции

Определение 5.1. Рассмотрим набор различных точек $\{x_i\}_{i=0}^n$, $x_i \in [a, b]$. Пусть $\{y_i\}_{i=0}^n$ — значения некоторой функции $f : [a, b] \rightarrow \mathbb{R}$ в этих точках: $y_i = f(x_i)$. Рассмотрим также набор линейно независимых базисных функций $\varphi_i : [a, b] \rightarrow \mathbb{R}$, $i = \overline{0, n}$. *Задача интерполяции* заключается в нахождении функции

$$\varphi = \sum_{i=0}^n \alpha_i \varphi_i, \quad \alpha_i \in \mathbb{R},$$

такой, что

$$\varphi(x_i) = y_i \quad \forall i = \overline{0, n}. \quad (5.1)$$

Функция f называется *интерполируемой функцией*, φ — *интерполирующей функцией*, $\{x_i\}$ — *узлами интерполяции*, точки декартовой плоскости (x_i, y_i) — *точками интерполяции*.

Замечание 5.1. Строго говоря, мы дали здесь определение *линейной* задачи интерполяции. В общем случае интерполирующая функция φ может зависеть от коэффициентов α_i нелинейно.

Таким образом, задача интерполяции сводится к нахождению неизвестных коэффициентов $\{\alpha_i\}_{i=0}^n$ из условий (5.1). По определению φ это эквивалентно решению СЛАУ

$$\sum_{j=0}^n \alpha_j \varphi_j(x_i) = y_i, \quad i = \overline{0, n},$$

матричный вид которой запишем как

$$\Phi \alpha = y, \quad (5.2)$$

где $\alpha = (\alpha_0, \dots, \alpha_n)^T$; $y = (y_0, \dots, y_n)^T$;

$$\Phi = \begin{bmatrix} \varphi_0(x_0) & \varphi_1(x_0) & \dots & \varphi_n(x_0) \\ \varphi_0(x_1) & \varphi_1(x_1) & \dots & \varphi_n(x_1) \\ \dots & \dots & \dots & \dots \\ \varphi_0(x_n) & \varphi_1(x_n) & \dots & \varphi_n(x_n) \end{bmatrix}. \quad (5.3)$$

Очевидно, что решение задачи интерполяции существует и единственно тогда и только тогда, когда $\det \Phi \neq 0$. Система (5.2) имеет наиболее простой вид в случае, когда $\Phi = I$. Отсюда вытекает следующее определение.

Определение 5.2. Пусть $\{x_i\}_{i=0}^n$ — узлы интерполяции. Система функций $\{\varphi_i\}_{i=0}^n$ называется *фундаментальным базисом* для данного набора узлов, если

$$\varphi_i(x_j) = \delta_{ij} = \begin{cases} 1, & i = j; \\ 0, & i \neq j; \end{cases} \quad \forall i, j = \overline{0, n}.$$

Таким образом, если $\{\varphi_i\}$ — фундаментальный базис, то задача интерполяции решается очень просто:

$$\varphi = \sum_{i=0}^n y_i \varphi_i.$$

5.1.2. Полиномиальная интерполяция

Классическим способом аппроксимации (приближения) функций является интерполяция алгебраическими многочленами:

$$\varphi(x) = P_n(x) = \sum_{i=0}^n \alpha_i x^i, \quad \text{т. е.} \quad \varphi_i(x) = x^i. \quad (5.4)$$

В этом случае интерполирующую функцию φ называют *интерполяционным многочленом*.

Пусть заданы узлы интерполяции. Рассмотрим для базиса из (5.4) матрицу интерполяции (5.3):

$$\Phi = \begin{bmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & \dots & x_2^n \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{bmatrix} = V(x_0, \dots, x_n). \quad (5.5)$$

Матрица такого вида называется *матрицей Вандермонда*.

Лемма 5.1. *Определитель матрицы Вандермонда может быть вычислен по формуле*

$$|V(x_0, \dots, x_n)| = \prod_{0 \leq j < i \leq n} (x_i - x_j). \quad (5.6)$$

Доказательство. Рассмотрим

$$|V(x_0, \dots, x_n)| = \begin{vmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & \dots & x_2^n \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{vmatrix}.$$

Преобразуем этот определитель: вычтем из j -го столбца $(j - 1)$ -й, умноженный на x_0 , для j от $n + 1$ до 2. В результате получим

$$|V(x_0, \dots, x_n)| = \begin{vmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & x_1 - x_0 & x_1^2 - x_1 x_0 & \dots & x_1^n - x_1^{n-1} x_0 \\ 1 & x_2 - x_0 & x_2^2 - x_2 x_0 & \dots & x_2^n - x_2^{n-1} x_0 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n - x_0 & x_n^2 - x_n x_0 & \dots & x_n^n - x_n^{n-1} x_0 \end{vmatrix} =$$

$$= (x_1 - x_0)(x_2 - x_0) \dots (x_n - x_0) \underbrace{\begin{vmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} \end{vmatrix}}_{|V(x_1, \dots, x_n)|}.$$

Преобразуя аналогичным образом по рекурсии $|V(x_1, \dots, x_n)|$, $|V(x_2, \dots, x_n)|$ и т. д., окончательно получим тождество (5.6). \square

Следствие 5.1. *Если все узлы интерполяции $\{x_i\}$ различны, то при любых значениях $\{y_i\}$ интерполяционный многочлен существует и единствен.*

Следствие 5.2. *Любой многочлен степени n однозначно определяется своими значениями в $n + 1$ различных точках.*

▷₁ Докажите данные следствия.

5.1.3. Интерполяционный многочлен в форме Лагранжа

Построим фундаментальный базис многочленов для узлов x_0, \dots, x_n . Пусть Λ_i — i -й многочлен из этого базиса. По определению

$$\Lambda_i(x_j) = 0 \quad \forall j \neq i,$$

т. е.

$$\Lambda_i(x) = C_i \prod_{j \neq i} (x - x_j), \quad C_i \in \mathbb{R}.$$

Здесь и далее используется обозначение

$$\prod_{j \neq i} (\dots) = \prod_{\substack{j=0 \\ j \neq i}}^n (\dots).$$

Неизвестную константу C_i найдем из оставшегося условия

$$\Lambda_i(x_i) = 1 \quad \Rightarrow \quad C_i = \left(\prod_{j \neq i} (x_i - x_j) \right)^{-1},$$

откуда

$$\Lambda_i(x) = \prod_{j \neq i} \frac{x - x_j}{x_i - x_j}. \quad (5.7)$$

Таким образом, получим знаменитую формулу *интерполяционного многочлена в форме Лагранжа*:

$$P_n(x) = \sum_{i=0}^n y_i \Lambda_i(x) = \sum_{i=0}^n y_i \prod_{j \neq i} \frac{x - x_j}{x_i - x_j}. \quad (5.8)$$

Для многочленов Λ_i существует альтернативная форма записи. Рассмотрим многочлен

$$\omega_{n+1}(x) = (x - x_0)(x - x_1) \dots (x - x_n). \quad (5.9)$$

Тогда числитель в (5.7) можно записать как

$$\frac{\omega_{n+1}(x)}{x - x_i},$$

а знаменатель как

$$\omega'_{n+1}(x_i).$$

Отсюда

$$\Lambda_i(x) = \frac{\omega_{n+1}(x)}{(x - x_i)\omega'_{n+1}(x_i)}. \quad (5.10)$$

5.1.4. Барицентрическая интерполяционная формула

Интерполяционная формула Лагранжа (5.8) становится громоздкой и неудобной для вычислений при больших n . Поэтому долгое время считалось, что эта формула малоприспособна для практических вычислений. Но оказывается, что при правильном способе записи ситуация кардинально изменяется.

Запишем многочлен Лагранжа с использованием формулы (5.10):

$$P_n(x) = \sum_{i=0}^n y_i \frac{\omega_{n+1}(x)}{(x - x_i)\omega'_{n+1}(x_i)}. \quad (5.11)$$

Чтобы избежать громоздких умножений в каждом слагаемом формулы (5.11), достаточно просто вынести общий множитель $\omega_{n+1}(x)$:

$$P_n(x) = \omega_{n+1}(x) \sum_{i=0}^n \frac{y_i}{(x - x_i)\omega'_{n+1}(x_i)}.$$

Эту формулу удобно записать в виде

$$P_n(x) = \omega_{n+1}(x) \sum_{i=0}^n y_i \frac{v_i}{x - x_i}, \quad (5.12)$$

где

$$v_i = \frac{1}{\omega'_{n+1}(x_i)} = \frac{1}{\prod_{j \neq i} (x_i - x_j)}, \quad i = \overline{0, n}. \quad (5.13)$$

Коэффициенты v_i будем называть *весовыми*. Полученная форма записи интерполяционного многочлена уже намного удобнее традиционной

формулы Лагранжа, но ее можно еще улучшить. Для этого построим по формуле (5.12) интерполяционный многочлен для $f(x) = 1$. Этот многочлен, как нетрудно заметить, тождественно равен единице при любом n :

$$\omega_{n+1}(x) \sum_{i=0}^n \frac{v_i}{x - x_i} = 1.$$

Выражая отсюда $\omega_{n+1}(x)$ и подставляя в (5.12), получим замечательную формулу

$$P_n(x) = \frac{\sum_{i=0}^n y_i \frac{v_i}{x - x_i}}{\sum_{i=0}^n \frac{v_i}{x - x_i}}, \quad (5.14)$$

которая называется *барицентрической интерполяционной формулой*. Обращаем внимание, что несмотря на то, что правая часть (5.14) выглядит как рациональная функция, по построению это не что иное, как интерполяционный многочлен.

Легко убедиться, что барицентрическая интерполяционная формула обладает следующими свойствами:

- для ее построения, т. е. для вычисления коэффициентов $\{v_i\}_{i=0}^n$, требуется $O(n^2)$ операций. Однако в п. 5.3 мы увидим, что при специальном выборе узлов интерполяции количество операций составляет всего $O(n)$;
- для вычисления $P_n(x)$ (при известных весовых коэффициентах) по этой формуле требуется всего $O(n)$ операций;
- формула может быть легко обновлена при добавлении дополнительного узла интерполяции x_{n+1} : каждое v_i необходимо домножить на $(x_i - x_{n+1})^{-1}$, после чего по определению вычислить v_{n+1} ;
- весовые коэффициенты не зависят от $\{y_i\}$, т. е. однажды построенная формула для заданного набора узлов может легко применяться для интерполяции различных функций f ;
- весовые коэффициенты v_i можно спокойно домножать на любую константу. Это большой плюс с точки зрения вычислительной устойчивости.

Помимо указанных свойств, барицентрическая формула обладает выдающейся вычислительной устойчивостью (при условии отдельной обработки случая, когда знаменатель в формуле (5.14) обращается в нуль).

5.2. Интерполяционный многочлен в форме Ньютона

5.2.1. Построение

Рассмотрим еще раз СЛАУ (5.2), решение которой является решением задачи интерполяции в общем случае:

$$\Phi \alpha = \begin{bmatrix} \varphi_0(x_0) & \varphi_1(x_0) & \dots & \varphi_n(x_0) \\ \varphi_0(x_1) & \varphi_1(x_1) & \dots & \varphi_n(x_1) \\ \dots & \dots & \dots & \dots \\ \varphi_0(x_n) & \varphi_1(x_n) & \dots & \varphi_n(x_n) \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

При построении формулы Лагранжа мы выбирали базис $\{\varphi_i\}$ так, чтобы матрица Φ была единичной. Теперь же построим базис, в котором эта матрица будет *нижнетреугольной*. Такой базис должен удовлетворять очевидным условиям

$$\varphi_i(x_j) = 0 \quad \forall j = \overline{0, i-1},$$

откуда с точностью до постоянного множителя получим $\varphi_i(x) = \omega_i(x)$, где

$$\omega_i(x) = (x - x_0)(x - x_1) \dots (x - x_{i-1}) = \prod_{j=0}^{i-1} (x - x_j) \quad (5.15)$$

(сравните с формулой (5.9)). Здесь при $i = 0$ полагаем $\omega_0(x) = 1$.

Определение 5.3. Интерполяционный многочлен вида

$$P_n(x) = \sum_{i=0}^n \alpha_i \omega_i(x) = \alpha_0 + \alpha_1(x - x_0) + \alpha_2(x - x_0)(x - x_1) + \dots + \alpha_n(x - x_0)(x - x_1) \dots (x - x_{n-1}) \quad (5.16)$$

называется *интерполяционным многочленом в форме Ньютона*. Правило вычисления коэффициентов α_i будет рассмотрено ниже.

Следует понимать, что формулы (5.16), (5.14) и (5.8) суть разные формы записи одного и того же многочлена.

По определению интерполяции коэффициенты интерполяционного многочлена в форме Ньютона удовлетворяют СЛАУ

$$\begin{bmatrix} \omega_0(x_0) & 0 & 0 & \dots & 0 \\ \omega_0(x_1) & \omega_1(x_1) & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ \omega_0(x_n) & \omega_1(x_n) & \omega_2(x_n) & \dots & \omega_n(x_n) \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}. \quad (5.17)$$

▷₁ При каких условиях СЛАУ с такой матрицей имеет единственное решение?

Следовательно, α_i можно найти, напрямую решив (5.17) путем обратной подстановки, которая приводит к следующим рекуррентным соотношениям:

$$\alpha_i = \left(y_i - \sum_{j=0}^{i-1} \alpha_j \omega_j(x_i) \right) / \omega_i(x_i), \quad i = \overline{0, n}. \quad (5.18)$$

Однако использование формул (5.18) не является общепринятым способом вычисления α_i , так как существует более эффективный способ, основанный на понятии разделенных разностей (см. далее).

По построению интерполяционный многочлен в форме Ньютона обладает рядом полезных свойств. Во-первых, базисные функции ω_i зависят только от узлов x_0, \dots, x_{i-1} , а коэффициенты α_i — только от x_0, \dots, x_i . Это позволяет легко «обновлять» формулу многочлена при добавлении дополнительных узлов интерполяции. Другими словами, если известен многочлен P_n (5.16), интерполирующий f по узлам x_0, \dots, x_n , то интерполяционный многочлен по узлам x_0, \dots, x_{n+1} может быть найден по формуле

$$P_{n+1}(x) = P_n(x) + \alpha_{n+1} \omega_{n+1}(x). \quad (5.19)$$

Во-вторых, форма (5.16) более удобна для вычислений, чем классическая форма Лагранжа: последовательно вынося за скобки общие множители $(x - x_i)$, можно вычислить $P_n(x)$ по так называемой схеме Горнера

$$P_n(x) = \alpha_0 + (x - x_0) \left(\alpha_1 + (x - x_1) (\alpha_2 + \dots) \right). \quad (5.20)$$

▷₂ Подсчитайте количество арифметических операций, нужных для вычисления $P_n(x)$ по формуле (5.20).

5.2.2. Разделенные разности

Пусть $\{y_i\}_{i=0}^n$ являются значениями некоторой функции f в точках $\{x_i\}_{i=0}^n$ соответственно. Тогда каждый из коэффициентов α_i , определяемых по формулам (5.18), можно рассматривать как выражение, зависящее от функции f и узлов x_0, \dots, x_i . Обозначим (пока формально) это выражение как

$$\alpha_i = f[x_0, \dots, x_i]. \quad (5.21)$$

Лемма 5.2. Пусть $\{x_0, x_1, \dots\}$ — последовательность различных узлов интерполяции, P_k — интерполяционный многочлен для функции f по узлам $\{x_0, \dots, x_k\}$, и P_k^+ — аналогичный многочлен, но построенный по узлам $\{x_1, \dots, x_{k+1}\}$. Тогда

$$P_{k+1}(x) = \frac{(x - x_0)P_k^+(x) - (x - x_{k+1})P_k(x)}{x_{k+1} - x_0}, \quad \forall k \geq 0. \quad (5.22)$$

▷₃ Докажите лемму простой проверкой условий $P_{k+1}(x_j) = f(x_j)$, $j = \overline{0, k+1}$.

Для $k = i - 1$ запишем интерполяционные многочлены P_k и P_k^+ из леммы 5.2 в форме Ньютона, используя обозначение (5.21):

$$\begin{aligned} P_{i-1}(x) &= \sum_{j=0}^{i-1} \alpha_j \omega_j(x) = \sum_{j=0}^{i-1} f[x_0, \dots, x_j] \omega_j(x), \\ P_{i-1}^+(x) &= \sum_{j=0}^{i-1} \alpha_j^+ \omega_j^+(x) = \sum_{j=0}^{i-1} f[x_1, \dots, x_{j+1}] \omega_j^+(x). \end{aligned}$$

Здесь $\omega_j^+(x) = (x - x_1)(x - x_2) \dots (x - x_j)$. Тогда (5.22) примет вид

$$\begin{aligned} P_i(x) &= \sum_{j=0}^i f[x_0, \dots, x_j] \omega_j(x) = \\ &= \frac{(x - x_0) \sum_{j=0}^{i-1} f[x_1, \dots, x_{j+1}] \omega_j^+(x) - (x - x_i) \sum_{j=0}^{i-1} f[x_0, \dots, x_j] \omega_j(x)}{x_i - x_0}. \end{aligned}$$

Приравнивая коэффициенты при старшей степени в обеих частях, получим важное тождество

$$f[x_0, \dots, x_i] = \frac{f[x_1, \dots, x_i] - f[x_0, \dots, x_{i-1}]}{x_i - x_0}.$$

Отсюда вытекает следующее определение.

Определение 5.4. Разделенной разностью порядка i для функции f по различным узлам $\{x_j\}_{j=0}^i$ называется выражение $f[x_0, \dots, x_i]$, определяемое по рекуррентным соотношениям

$$f[x_0, \dots, x_i] = \frac{f[x_1, \dots, x_i] - f[x_0, \dots, x_{i-1}]}{x_i - x_0}, \quad \forall i \geq 1; \quad (5.23a)$$

$$f[x_j] = f(x_j) \quad \forall j. \quad (5.23б)$$

Таким образом, коэффициенты (5.21) интерполяционного многочлена в форме Ньютона являются разделенными разностями и традиционно вычисляются по определению (5.23).

▷₄ Вычислите общий вид коэффициентов α_0 , α_1 и α_2 сначала по формуле (5.18), затем по формуле (5.23).

▷₅ Сравните вычислительную сложность вычисления α_i по формулам (5.18) и (5.23).

▷₆ Докажите, что значение разделенной разности не зависит от порядка расположения ее аргументов.

5.2.3. Алгоритм вычисления разделенных разностей

Коэффициенты интерполяционного многочлена (5.16) удобно вычислять по определению разделенных разностей (5.23) путем построения треугольной таблицы следующего вида:

x_0	$f[x_0]$			
		$f[x_0, x_1]$		
x_1	$f[x_1]$		$f[x_0, x_1, x_2]$	
		$f[x_1, x_2]$	\dots	
x_2	$f[x_2]$		$f[x_1, x_2, x_3]$	$\dots f[x_0, \dots, x_{n-1}]$
	\cdot	$f[x_2, x_3]$	\cdot	$\dots f[x_0, \dots, x_n]$
	\cdot	\cdot	\cdot	$\dots f[x_1, \dots, x_n]$
	\cdot	\cdot	\cdot	\dots
x_{n-1}	$f[x_{n-1}]$	\cdot	$f[x_{n-2}, x_{n-1}, x_n]$	
		$f[x_{n-1}, x_n]$		
x_n	$f[x_n]$			

Вычисления осуществляются следующим образом: сначала записываются первые два столбца таблицы, затем по определению вычисляется третий (разделенные разности первого порядка), затем четвертый и

т. д. При этом каждый новый элемент таблицы вычисляется с использованием двух ближайших к нему элементов предыдущего столбца и соответствующих им значений x_i .

▷₇ Каким минимальным количеством ячеек памяти можно обойтись при вычислении разделенных разностей, необходимых для построения интерполяционного многочлена в форме Ньютона? Постройте соответствующий алгоритм.

▷₈ Сравните трудоемкость построения и вычисления интерполяционного многочлена в форме Ньютона с барицентрической интерполяционной формулой (5.14).

5.3. Остаток алгебраической интерполяции

5.3.1. Остаток алгебраической интерполяции

Определение 5.5. C -нормой для непрерывной на отрезке $[a, b]$ функции f называется величина $\|f\|_C$, вычисляемая по формуле

$$\|f\|_C = \|f\|_{C[a,b]} = \max_{a \leq x \leq b} |f(x)|.$$

Определение 5.6. Пусть $f \in C[a, b]$ — интерполируемая функция, P_n — интерполяционный многочлен для f по узлам $\{x_i \in [a, b]\}_{i=0}^n$. *Остатком интерполирования* называют функцию

$$r_n = f - P_n. \quad (5.24)$$

Погрешностью интерполирования назовем C -норму остатка:

$$\varepsilon_n = \|r_n\|_C = \max_{x \in [a,b]} |r_n(x)|. \quad (5.25)$$

Теорема 5.1. *Остаток интерполирования имеет вид*

$$r_n(x) = f[x_0, \dots, x_n, x] \omega_{n+1}(x) \quad \forall x \in \mathbb{R}. \quad (5.26)$$

Доказательство. Пусть P_n — интерполяционный многочлен по узлам x_0, \dots, x_n для f . Рассмотрим произвольную точку $\xi \neq x_i, i = \overline{0, n}$. Построим для f интерполяционный многочлен P_{n+1} по узлам $\{x_0, \dots, x_n, \xi\}$. По формуле (5.19) имеем

$$P_{n+1}(\xi) = f(\xi) = P_n(\xi) + f[x_0, \dots, x_n, \xi] \omega_{n+1}(\xi).$$

Меняя в этой формуле обозначение ξ на x , получим (5.26). □

Замечание 5.2. Если точка x совпадает с одним из узлов интерполяции x_i , то в формуле (5.26) возникнет разделенная разность с двумя одинаковыми узлами, для вычисления которой нужно обобщить определение 5.4. Такое обобщение можно найти в большинстве пособий по численному анализу, например в [2]. Для строгого же доказательства теоремы 5.1 нам достаточно принять на веру, что такая разделенная разность будет ограниченной величиной: тогда (5.26) в соответствии с реальностью даст $r_n(x) = 0$.

Чтобы получить еще одно полезное выражение для остатка, необходимо вспомнить следующую теорему из курса анализа.

Теорема 5.2 (Ролля). Если функция g непрерывна на $[a, b]$, дифференцируема на (a, b) , и $g(a) = g(b)$, то существует по крайней мере одна такая точка $\xi \in (a, b)$, что $g'(\xi) = 0$.

С помощью этой теоремы можно установить связь между разделенными разностями и производными функции f .

Лемма 5.3. Пусть $\underline{x} = \min\{x_i\}_{i=0}^n$, $\bar{x} = \max\{x_i\}_{i=0}^n$ и $f \in C^n[\underline{x}, \bar{x}]$. Тогда $\exists \xi \in (\underline{x}, \bar{x})$ такое, что

$$f[x_0, \dots, x_n] = \frac{f^{(n)}(\xi)}{n!}. \quad (5.27)$$

Доказательство. Построим для f интерполяционный многочлен P_n в форме Ньютона по узлам x_0, \dots, x_n и рассмотрим остаток $r_n = f - P_n$. По построению имеем

$$r_n(x_0) = r_n(x_1) = \dots = r_n(x_n) = 0,$$

значит, по теореме Ролля между \underline{x} и \bar{x} существует как минимум n точек, в которых r'_n обращается в нуль.

Продолжая аналогичные рассуждения для r''_n и т. д., получим, что существует $\xi \in (\underline{x}, \bar{x})$, такое, что

$$r_n^{(n)}(\xi) = f^{(n)}(\xi) - P_n^{(n)}(\xi) = f^{(n)}(\xi) - n! f[x_0, \dots, x_n] = 0. \quad \square$$

Теорема 5.3. Пусть $f \in C^{(n+1)}[a, b]$, P_n — интерполяционный многочлен для f по узлам $\{x_i\}_{i=0}^n \subset [a, b]$. Тогда $\forall x \in [a, b] \exists \xi \in (\underline{x}, \bar{x})$, $(\underline{x} = \min\{x_0, \dots, x_n, x\}, \bar{x} = \max\{x_0, \dots, x_n, x\})$, такое, что

$$r_n(x) = f(x) - P_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{n+1}(x). \quad (5.28)$$

Доказательство. Данная теорема представляет собой тривиальное следствие теоремы 5.1 и леммы 5.3. \square

Формула (5.28) является основным инструментом при оценке погрешности интерполяции. Из нее, в частности, сразу получим

$$\varepsilon_n = \|r_n\| \leq \frac{\|f^{(n+1)}\|}{(n+1)!} \|\omega_{n+1}\|, \quad (5.29)$$

где $\|\cdot\| = \|\cdot\|_{C[a,b]}$ (см. (5.25)).

5.3.2. Задача выбора оптимальных узлов интерполирования

Пусть P_n — интерполяционный многочлен для f по узлам $\{x_i\}_{i=0}^n$ из отрезка $[a, b]$. Согласно теореме 5.3, точность приближения оценивается по формуле

$$\|f - P_n\| \leq \frac{\|f^{(n+1)}\|}{(n+1)!} \|\omega_{n+1}\|, \quad (5.30)$$

где

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

Здесь и далее в этом подпункте $\|\cdot\| = \|\cdot\|_C$.

Часто требуется приблизить заранее неизвестную функцию интерполяционным многочленом как можно точнее, т. е. сделать правую часть (5.30) как можно меньше. И если в общем случае о величине $\|f^{(n+1)}\|$ ничего не известно, то $\|\omega_{n+1}\|$ можно минимизировать за счет выбора узлов $\{x_i\}$. Задачу выбора оптимальных узлов интерполяции можно поставить так: найти многочлен ω_{n+1} со *старшим коэффициентом, равным единице*, и минимально возможной нормой $\|\omega_{n+1}\|$ на отрезке интерполяции $[a, b]$. Такой многочлен называется *многочленом, наименее отклоняющимся от нуля* на $[a, b]$. Корни такого многочлена и будут оптимальными⁶ узлами интерполяции.

Оказывается, что решением поставленной задачи для отрезка $[-1, 1]$ являются многочлены Чебышева, которые мы рассматривали в разделе 2.8.2.

5.3.3. Оптимальные узлы интерполирования

Нули многочлена Чебышева (2.47) являются оптимальными узлами интерполирования на отрезке $[-1, 1]$, т. е. оптимальный многочлен ω_{n+1} определяется формулой

$$\omega_{n+1}(x) = 2^{-n} T_{n+1}(x). \quad (5.31)$$

Множитель 2^{-n} — это обратное значение старшего коэффициента многочлена T_{n+1} , в чем нетрудно убедиться из рекуррентного соотноше-

⁶Необходимо понимать, что «оптимальные» в данном случае означает «универсальные», не зависящие от вида интерполируемой функции. Для каждой конкретной функции f узлы интерполяции, дающие минимальную погрешность, вообще говоря, заранее не известны.

ния (2.48). Нули этого многочлена легко находятся из уравнения

$$T_{n+1}(x) = \cos((n+1) \arccos x) = 0$$

и равны

$$x_i = \cos \frac{\pi(2i+1)}{2n+2}, \quad i = \overline{0, n}. \quad (5.32)$$

▷₁ Как надо изменить формулу (5.32), чтобы узлы интерполяции расположились в порядке возрастания?

Это и есть оптимальные узлы интерполирования, называемые также *чебышевскими*. При таком выборе узлов из (5.31) получим

$$\|\omega_{n+1}\| = 2^{-n},$$

т. е. оценка погрешности (5.33) для интерполяционного многочлена по чебышевским узлам на отрезке $[-1, 1]$ имеет вид

$$\|r_n\| \leq \frac{\|f^{(n+1)}\|}{2^n(n+1)!}. \quad (5.33)$$

Весовые коэффициенты барицентрической формулы

Использование чебышевских узлов интерполяции, помимо прочего, существенно повышает эффективность барицентрической интерполяционной формулы (5.14). Напомним, что для ее использования необходимо вычислять весовые коэффициенты по формуле

$$v_i = \frac{1}{\prod_{j \neq i} (x_i - x_j)}, \quad i = \overline{0, n}.$$

Если в этой формуле использовать чебышевские узлы (5.32), то с учетом (5.31), после сокращения множителей, не зависящих от i , она принимает простой вид

$$v_i = (-1)^i \sin \frac{2i+1}{2n+2} \pi. \quad (5.34)$$

▷₂ Выведите эту формулу

Еще более простой вид имеют коэффициенты для так называемых *чебышевских узлов второго рода*

$$x_i = \cos \frac{i\pi}{n}, \quad i = \overline{0, n}. \quad (5.35)$$

Соответствующие весовые коэффициенты для барицентрической формулы определяются формулой

$$v_i = (-1)^i \delta_i, \quad \delta_i = \begin{cases} 1/2, & \text{если } i = 0 \text{ или } i = n, \\ 1, & \text{в противном случае.} \end{cases} \quad (5.36)$$

Как видим, для узлов такого типа барицентрическая интерполяционная формула практически не требует никаких предварительных вычислений.

Чебышевские узлы на отрезке $[a, b]$

Чтобы применить полученные результаты для интерполяции на произвольном отрезке $[a, b]$, нужно воспользоваться заменой переменной:

$$x = \frac{a+b}{2} + \frac{b-a}{2}t, \quad t \in [-1, 1], \quad x \in [a, b].$$

Тогда многочлен Чебышева, смасштабированный на $[a, b]$, имеет вид

$$\begin{aligned} \hat{T}_{n+1}(x) &= T_{n+1}(t) = 2t T_n(t) - T_{n-1}(t) = 2t \hat{T}_n(x) - \hat{T}_{n-1}(x) = \\ &= 2 \frac{2x - a - b}{b - a} \hat{T}_n(x) - \hat{T}_{n-1}(x). \end{aligned} \quad (5.37)$$

Кроме этого, имеют место очевидные соотношения

$$\hat{T}_0(x) = 1, \quad \hat{T}_1(x) = \frac{2x - b - a}{b - a}.$$

Следовательно, согласно (5.37), старший коэффициент многочлена $\hat{T}_n(x)$ при всех $n \geq 1$ равен

$$\hat{a}_n = \frac{2}{b-a} \left(\frac{4}{b-a} \right)^{n-1} = \frac{1}{2} \left(\frac{4}{b-a} \right)^n,$$

откуда получим

$$\omega_{n+1}(x) = 2 \left(\frac{b-a}{4} \right)^{n+1} T_{n+1} \left(\frac{2x - b - a}{b - a} \right). \quad (5.38)$$

Корни этого многочлена (*оптимальные узлы интерполяции на отрезке $[a, b]$*), получаются масштабированием узлов (5.32):

$$x_i = \frac{a+b}{2} + \frac{b-a}{2} \cos \frac{\pi(2i+1)}{2n+2}, \quad i = \overline{0, n}. \quad (5.39)$$

Кроме этого, из (5.38) следует, что при выборе чебышевских узлов интерполирования (5.39) имеем

$$\|\omega_{n+1}\| = 2 \left(\frac{b-a}{4} \right)^{n+1}.$$

Это равенство используется при оценке погрешности интерполирования по формуле (5.30):

$$\|r_n\| \leq 2 \left(\frac{b-a}{4} \right)^{n+1} \frac{\|f^{(n+1)}\|}{(n+1)!}. \quad (5.40)$$

5.3.4. О сходимости интерполяционного процесса

Центральный вопрос теории приближения функций — сходимость интерполяционного процесса. Говоря простым языком, необходимо определить, будет ли интерполяционный многочлен стремиться к интерполируемой функции при неограниченном увеличении количества узлов интерполяции. На языке формул это требование записывается как

$$\|f - P_n\|_C \xrightarrow{n \rightarrow \infty} 0.$$

Здесь P_n — интерполяционный многочлен, построенный по набору узлов

$$X_n = \{x_0^{(n)}, x_1^{(n)}, \dots, x_n^{(n)}\}.$$

Таким образом, последовательность $\{P_n\}$, т. е. сам интерполяционный процесс, определяется таблицей узлов интерполяции

$$\begin{aligned} X_0 &= \{x_0^{(0)}\}, \\ X_1 &= \{x_0^{(1)}, x_1^{(1)}\}, \\ &\dots\dots\dots \\ X_n &= \{x_0^{(n)}, x_1^{(n)}, \dots, x_n^{(n)}\}, \\ &\dots\dots\dots \end{aligned} \quad (5.41)$$

Основные факты о сходимости полиномиальной интерполяции

Теорема 5.4. *Для любой $f \in C[a, b]$ существует такая последовательность сеток, для которой интерполяционный процесс равномерно сходится к f .*

Эта теорема носит формальный характер, так как построение нужной последовательности сеток в общем случае представляет серьезную проблему. Следующий результат говорит о том, что не существует «универсальной» последовательности сеток, хорошей для всех непрерывных функций.

Теорема 5.5. *Для любой последовательности сеток вида (5.41) существует такая $f \in C[a, b]$, для которой интерполяционный процесс не сходится равномерно к f .*

Таким образом, класс $C[a, b]$ слишком широк. Поэтому рассмотрим результаты о сходимости интерполяционного процесса для более узкого класса функций.

Определение 5.7. Функция f называется целой, если существует ее разложение в степенной ряд вида

$$f(x) = \sum_{i=0}^{\infty} c_i (x - x_0)^i,$$

которое сходится при любом x .

Типичные примеры целых функций: многочлены, \exp , \sin , \cos и их линейные комбинации.

Теорема 5.6. *Если функция f целая, то интерполяционный процесс по любой последовательности сеток вида (5.41) равномерно сходится на $[a, b]$ к f .*

Тем не менее класс целых функций достаточно узок. Например, интерполяция по равноотстоящим узлам расходится для известной функции Рунге $f(x) = (1 + 25x^2)^{-1}$ и функции $g(x) = |x|$, обе на отрезке $[-1, 1]$. Для последней интерполяционный многочлен по равноотстоящим узлам степени $2n$ неограниченно растет в любой части отрезка $[-1, 1]$.

Кроме этого, интерполяция по равноотстоящим узлам обладает еще одним неприятным свойством: она плохо обусловлена, т. е. при достаточно большом количестве точек небольшие погрешности округления будут приводить к большим погрешностям. Значит, даже если для функции \cos интерполяция по равноотстоящим узлам сходится *теоретически* по теореме 5.6, на практике (в машинной арифметике) с ростом числа узлов интерполяционный процесс будет расходиться.

Более привлекательна интерполяция по чебышевским узлам.

Теорема 5.7. *Для любой абсолютно непрерывной на $[a, b]$ функции f интерполяционный процесс по чебышевским узлам равномерно сходится.*

Точное определение абсолютно непрерывной функции нам не нужно, достаточно знать, что любая функция с ограниченной производной на отрезке абсолютно непрерывна на нем. Помимо этого, интерполяция по чебышевским узлам хорошо обусловлена.

5.4. Сплайны

5.4.1. Определение

Рассмотрим набор узлов

$$\{x_i\}_{i=0}^n, \quad a = x_0 < x_1 < \dots < x_{n-1} < x_n = b,$$

и соответствующее ему разбиение отрезка $[a, b]$ на n частей:

$$\Delta = \{\Delta_i\}_{i=1}^n, \quad \Delta_i = [x_{i-1}, x_i]. \quad (5.42)$$

Пример такого разбиения для случая $n = 4$ изображен на рис. 5.2. На каждом отрезке Δ_i зададим многочлен s_i фиксированной степени $m \geq 1$ и рассмотрим полученный набор многочленов как единую кусочно-непрерывную функцию s , определенную на всем отрезке $[a, b]$:

$$s(x) = s_i(x), \quad \text{если } x \in \Delta_i.$$

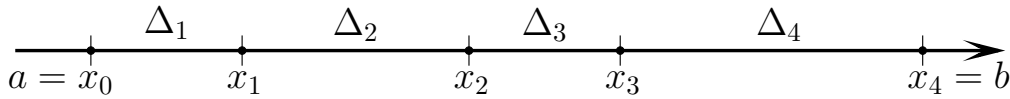


Рис. 5.2. Пример разбиения отрезка $[a, b]$

Для того чтобы это определение было однозначным при $x = x_i$, будем считать, что $s_i(x_i) = s_{i+1}(x_i)$ для всех $i = \overline{1, n-1}$, т. е. что функция s непрерывна. Такая кусочная функция и называется полиномиальным сплайном. Сформулируем это определение более строго.

Определение 5.8. Пусть Δ — некоторое разбиение отрезка $[a, b]$, определяемое формулой (5.42), \mathbb{P}_m — множество всех многочленов степени не выше m . Для некоторого $m \geq 1$ рассмотрим функцию s , определенную на $[a, b]$ и обладающую следующими свойствами:

- 1) $s(x)|_{x \in \Delta_i} = s_i(x), \quad s_i \in \mathbb{P}_m$;
- 2) функция s имеет $m - 1$ непрерывных производных:

$$s \in C^{m-1}[a, b],$$

что эквивалентно условиям

$$\begin{aligned} s_i^{(j)}(x_i - 0) &= s_{i+1}^{(j)}(x_i + 0), \\ i &= \overline{1, n-1}, \quad j = \overline{0, m-1}. \end{aligned} \quad (5.43)$$

Такая функция s называется *полиномиальным сплайном порядка (степени) m* . Множество всех сплайнов степени m на разбиении Δ будем обозначать S_{Δ}^m .

5.4.2. Интерполяционные сплайны

Определение 5.9. Сплайн $s \in S_{\Delta}^m$ называется интерполяционным для функции f , если

$$s(x_i) = f(x_i) = y_i \quad \forall i = \overline{0, n}. \quad (5.44)$$

Интерполяционный сплайн $s \in S_{\Delta}^1$ представляет собой кусочно-линейную функцию (график — ломаная линия, рис. 5.3), построенную по точкам $\{(x_i, y_i)\}_{i=0}^n$.

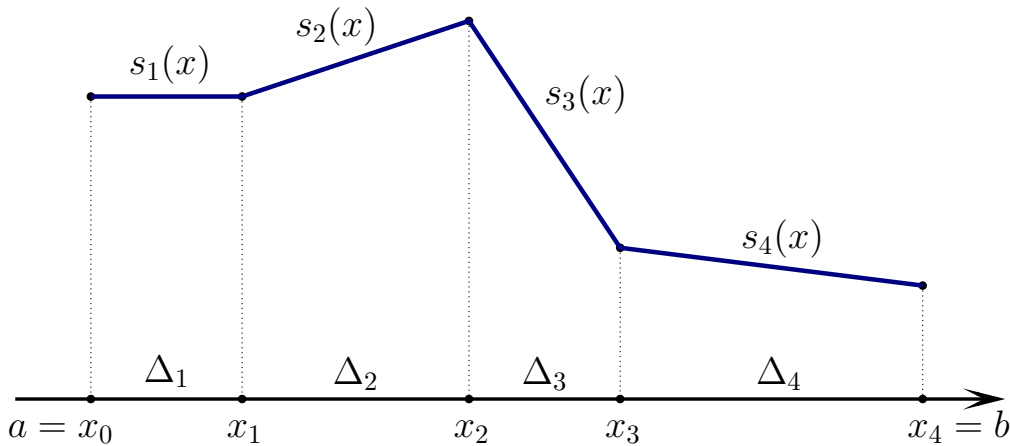


Рис. 5.3. Интерполяционный линейный сплайн

Большой интерес, конечно, представляют сплайны высших степеней. Наиболее распространенными являются кубические интерполяционные сплайны ($m = 3$, рис. 5.4).

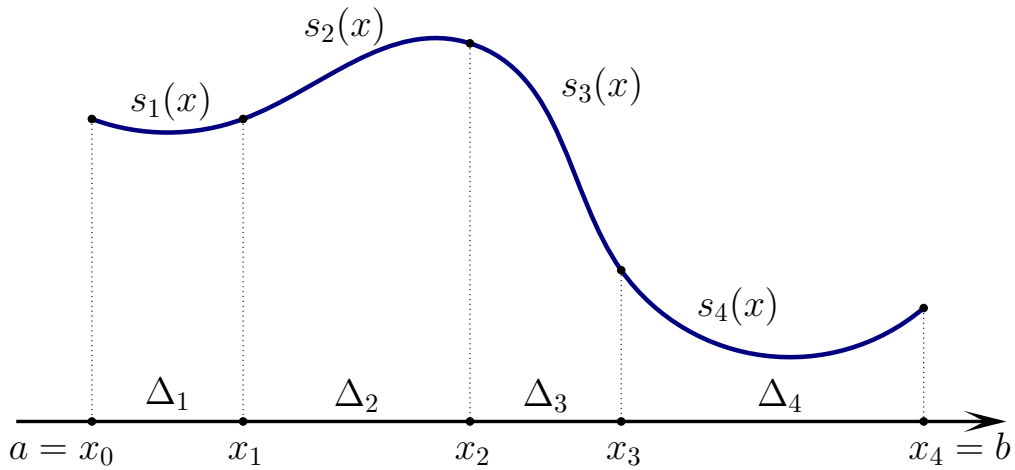


Рис. 5.4. Интерполяционный кубический сплайн

Построение кубического интерполяционного сплайна

Рассмотрим задачу нахождения $s \in S_{\Delta}^3$, удовлетворяющего условиям (5.44). Каждый «кусок» сплайна будем искать в виде

$$s_i(x) = \alpha_i + \beta_i(x - x_i) + \frac{\gamma_i}{2}(x - x_i)^2 + \frac{\delta_i}{6}(x - x_i)^3, \quad x \in \Delta_i = [x_{i-1}, x_i].$$

Такое представление в виде многочлена по степеням $(x - x_i)$ выбрано не случайно: при этом имеем

$$\begin{aligned} s'_i(x) &= \beta_i + \gamma_i(x - x_i) + \frac{\delta_i}{2}(x - x_i)^2, \\ s''_i(x) &= \gamma_i + \delta_i(x - x_i), \\ s'''_i(x) &= \delta_i \end{aligned}$$

и

$$s_i(x_i) = \alpha_i, \quad s'_i(x_i) = \beta_i, \quad s''_i(x_i) = \gamma_i.$$

Для построения всего сплайна нам нужно определить $4n$ неизвестных $\{\alpha_i, \beta_i, \gamma_i, \delta_i\}_{i=1}^n$. Обозначим

$$\alpha_0 = s_1(x_0), \quad \beta_0 = s'_1(x_0), \quad \gamma_0 = s''_1(x_0), \quad \delta_0 = s'''_1(x_0), \quad (5.45)$$

после чего с полным правом можем записать

$$s(x_i) = \alpha_i, \quad s'(x_i) = \beta_i, \quad s''(x_i) = \gamma_i, \quad \forall i = \overline{0, n}.$$

Выпишем уравнения, которым должны удовлетворять искомые коэффициенты. Обозначим $h_i = x_i - x_{i-1}$. Прежде всего используем условия интерполяции (5.44). Имеем $s_i(x_i) = y_i$, откуда

$$\alpha_i = y_i, \quad i = \overline{1, n}, \quad (5.46a)$$

и $s_i(x_{i-1}) = y_{i-1}$, или

$$\beta_i = \frac{y_i - y_{i-1}}{h_i} + \frac{\gamma_i}{2}h_i - \frac{\delta_i}{6}h_i^2, \quad i = \overline{1, n}. \quad (5.46б)$$

Отметим, что выполнение этих двух условий сразу гарантирует непрерывность s . Теперь используем условия непрерывности s' :

$$s'_i(x_{i-1}) = s'_{i-1}(x_{i-1}),$$

или

$$\beta_i - \gamma_i h_i + \frac{\delta_i}{2}h_i^2 = \beta_{i-1}, \quad i = \overline{2, n}. \quad (5.46в)$$

Наконец, требование непрерывности второй производной

$$s''_i(x_{i-1}) = s''_{i-1}(x_{i-1})$$

по аналогии дает

$$\delta_i = \frac{\gamma_i - \gamma_{i-1}}{h_i}, \quad i = \overline{2, n}. \quad (5.46г)$$

Важно: если в последних двух формулах положить $i = 1$, то мы получим просто определение величин β_0 и γ_0 согласно (5.45).

Формулы (5.46) дают $4n - 2$ уравнений для нахождения $4n$ неизвестных, поэтому без двух дополнительных условий сплайн однозначно не определить. Как правило, эти условия задаются на концах отрезка $[a, b]$ и поэтому называются *граничными*.

Величины γ_i называются *моментами* сплайна. Прежде чем рассматривать различные варианты задания дополнительных условий для определения интерполяционного кубического сплайна, заметим, что если известны значения всех моментов, то остальные коэффициенты сразу находятся по формулам (5.46a), (5.46б) и (5.46г) (в последней нужно взять еще $i = 1$). В частности, (5.46б) дает

$$\beta_i = \frac{y_i - y_{i-1}}{h_i} + \frac{2\gamma_i + \gamma_{i-1}}{6}h_i. \quad (5.46б')$$

Для нахождения $\{\gamma_i\}$ подставим (5.46б') и (5.46г) в (5.46в). После сдвига в нумерации получим

$$h_i \gamma_{i-1} + 2(h_i + h_{i+1})\gamma_i + h_{i+1}\gamma_{i+1} = 6 \left(\frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i} \right),$$

что можно записать как

$$c_i \gamma_{i-1} + 2\gamma_i + e_i \gamma_{i+1} = b_i, \quad i = \overline{1, n-1}, \quad (5.47)$$

где

$$c_i = \frac{h_i}{h_i + h_{i+1}}, \quad e_i = \frac{h_{i+1}}{h_i + h_{i+1}}, \quad b_i = 6f[x_{i-1}, x_i, x_{i+1}].$$

Формулы (5.47) дают $(n-1)$ линейных уравнений для нахождения $n+1$ неизвестных $\{\gamma_i\}_{i=0}^n$.

Таким образом, **общая схема вычисления интерполяционного кубического сплайна** такова:

- 1) вычисляются моменты $\{\gamma_i\}_{i=0}^n$ как решение СЛАУ (5.47), дополненной двумя дополнительными уравнениями — граничными условиями;
- 2) находятся остальные коэффициенты по формулам (5.46а), (5.46б') и (5.46г).

Виды граничных условий

Естественные граничные условия. Удобные граничные условия $s''(a) = s''(b) = 0$, как и соответствующий интерполяционный сплайн, называются *естественными*. В этом случае имеем $\gamma_0 = \gamma_n = 0$, а остальные $\{\gamma_i\}_{i=1}^{n-1}$ легко найти из (5.47), которая имеет вид

$$\begin{bmatrix} 2 & e_1 & & & \\ c_2 & 2 & e_2 & & \\ & c_3 & 2 & e_3 & \\ & & \ddots & \ddots & \ddots \\ & & & c_{n-2} & 2 & e_{n-2} \\ & & & & c_{n-1} & 2 \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \gamma_3 \\ \vdots \\ \gamma_{n-2} \\ \gamma_{n-1} \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_{n-2} \\ b_{n-1} \end{bmatrix}.$$

Данную систему решают методом прогонки.

▷₁ Докажите, что метод прогонки в данном случае применим.

▷₂ Запишите вид аналогичной СЛАУ для граничных условий вида $s''(x_0) = f''(x_0)$, $s''(x_n) = f''(x_n)$.

Граничные условия на значение первой производной. Пусть в концах отрезка сплайн s должен удовлетворять дополнительным условиям

$$s'(a) = f'(a), \quad s'(b) = f'(b),$$

что равносильно $\beta_0 = f'(a)$, $\beta_n = f'(b)$. Рассмотрим (5.46в) при $i = 1$ с учетом (5.46б):

$$\gamma_1 h_1 - \frac{\gamma_1 - \gamma_0}{2} h_1 = \frac{y_1 - y_0}{h_1} + \frac{2\gamma_1 + \gamma_0}{6} h_1 - f'(a),$$

откуда получим

$$2\gamma_0 + \gamma_1 = \frac{6}{h_1} \left(\frac{y_1 - y_0}{h_1} - f'(a) \right).$$

Аналогично для $i = n - 1$ имеем

$$\gamma_{n-1} + 2\gamma_n = \frac{6}{h_n} \left(f'(b) - \frac{y_n - y_{n-1}}{h_n} \right).$$

▷₃ Запишите соответствующую СЛАУ для нахождения $\{\gamma_i\}_{i=0}^n$ и докажите, что к ней применим метод прогонки.

Другие граничные условия. На практике используется еще несколько видов граничных условий, например периодические:

$$s'(a) = s'(b), \quad s''(a) = s''(b)$$

или так называемые условия «отсутствия узлов», накладывающие требования непрерывности s''' в узлах x_1 и x_{n-1} .

▷₄ Почему в этом случае говорят об «отсутствии узлов»?

Интересен также вариант

$$s'(a) = f'(a), \quad s''(a) = f''(a).$$

В этом случае вообще не нужно решать никакой СЛАУ — все s_i вычисляются явно по очереди от 1 до n .

5.4.3. Экстремальное свойство кубического сплайна

Рассмотрим функционал $\Phi : C^2[a, b] \rightarrow \mathbb{R}$,

$$\Phi(f) = \int_a^b (f''(x))^2 dx.$$

Экстремальное свойство интерполяционного сплайна $s \in S_{\Delta}^3$ состоит в том, что он доставляет минимум Φ среди всех функций из $C^2[a, b]$, интерполирующих f и удовлетворяющих требуемым граничным условиям. Сформулируем это утверждение более строго для случая естественных граничных условий.

Теорема 5.8. Пусть $F_0 \subset C^2[a, b]$ — множество функций f , таких, что

$$f(x_i) = y_i, \quad i = \overline{0, n}, \quad \text{и} \quad f''(a) = f''(b) = 0.$$

Если кубический сплайн $s \in S_{\Delta}^3$ принадлежит F_0 , то

$$\Phi(s) \leq \Phi(f) \quad \forall f \in F_0.$$

Доказательство. Для произвольной $f \in F_0$ рассмотрим функцию $g = f - s$. По условию имеем $g \in F_0$ и

$$g(x_i) = 0 \quad \forall i = \overline{0, n}.$$

Расписывая $\Phi(f) = \Phi(s + g)$, получим

$$\Phi(f) = \Phi(s) + \Phi(g) + 2 \int_a^b s'' g'' dx.$$

Для доказательства теоремы теперь достаточно показать, что

$$\int_a^b s'' g'' dx \geq 0.$$

Воспользуемся интегрированием по частям, а также тем фактом, что s'''

на интервале (x_{i-1}, x_i) равна константе, которую обозначим δ_i :

$$\begin{aligned}
 \int_a^b s'' g'' dx &= \sum_{i=1}^n \int_{x_{i-1}}^{x_i} s'' g'' dx = \sum_{i=1}^n \left(s'' g' \Big|_{x_{i-1}}^{x_i} - \int_{x_{i-1}}^{x_i} s''' g' dx \right) = \\
 &= \underbrace{s''(b)g'(b)}_0 - \underbrace{s''(a)g'(a)}_0 - \sum_{i=1}^n \int_{x_{i-1}}^{x_i} s''' g' dx = - \sum_{i=1}^n \int_{x_{i-1}}^{x_i} s''' g' dx = \\
 &= - \sum_{i=1}^n \delta_i \int_{x_{i-1}}^{x_i} g' dx = - \sum_{i=1}^n \delta_i (g(x_i) - g(x_{i-1})) = 0. \quad \square
 \end{aligned}$$

Отметим, что в аналитической геометрии кривизна кривой — графика $y = f(x)$ — вычисляется по формуле

$$\frac{|f''(x)|}{(1 + f'(x)^2)^{3/2}},$$

так что величину $\Phi(f)$ можно характеризовать как приближенное значение для усредненной кривизны графика. Поэтому часто говорят (пусть и не совсем строго), что интерполяционный сплайн обладает минимальной кривизной среди всех интерполирующих функций с указанными граничными условиями.

5.4.4. Сходимость интерполяции кубическим сплайном

Приведем без доказательства следующий результат.

Теорема 5.9. *Рассмотрим $f \in C^4[a, b]$ и разбиение Δ отрезка $[a, b]$, такое, что $\max_i h_i = h$. Пусть $s \in S_\Delta^3$ — естественный интерполяционный сплайн для f . Тогда*

$$\|f - s\| \leq \frac{5}{384} h^4 \|f^{(4)}\|.$$

Здесь $\|\cdot\| = \|\cdot\|_{C[a,b]}$.

Таким образом, для достаточно гладких f погрешность интерполирования кубическим сплайном $\varepsilon(h)$ равна $O(h^4)$.

5.5. Приближение кривых

5.5.1. Интерполяция кривых на плоскости

Постановка задачи

Рассмотрим на декартовой плоскости кривую ℓ . На данной кривой зададим направление движения и $n + 1$ точек $p_i = (x_i, y_i)^T$, $i = \overline{0, n}$, расположенных последовательно в выбранном направлении (рис. 5.5).

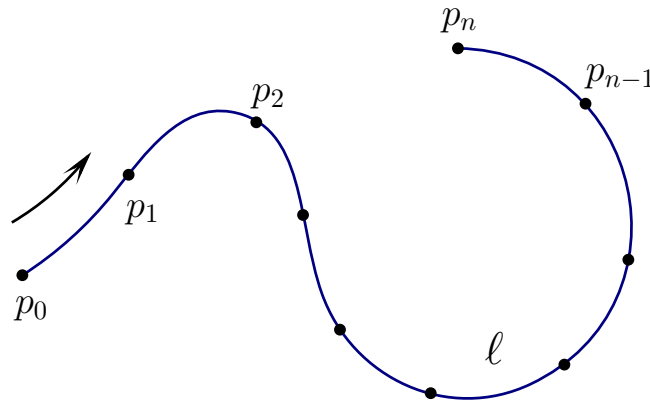


Рис. 5.5. Задача интерполяции кривой на плоскости

Определение 5.10. Задача интерполяции кривой ℓ по точкам $\{p_i\}$ заключается в построении такой параметрической функции $g : \mathbb{R} \rightarrow \mathbb{R}^2$,

$$g(t) = \begin{bmatrix} x(t) \\ y(t) \end{bmatrix}, \quad (5.48)$$

что

$$g(t_i) = p_i, \quad \forall i = \overline{0, n},$$

где $\{t_i\}$ — некоторая возрастающая последовательность значений параметра:

$$t_0 < t_1 < \dots < t_{n-1} < t_n,$$

которые будем называть *параметрическими узлами*.

В качестве параметрических узлов чаще всего используются следующие:

- *равномерные параметрические узлы* на отрезке $[0, n]$ или $[0, 1]$: $t_i = i$ и $t_i = i/n$ соответственно;

- *естественные параметрические узлы*, зависящие от расстояния между точками интерполяции p_i :

$$t_0 = 0, \quad t_{i+1} = t_i + \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2}, \quad i = \overline{1, n}.$$

Все эти определения, а также методы приближения плоских кривых, рассмотренные далее, могут быть естественным образом перенесены и на случай интерполяции кривой в пространстве размерности три и более.

5.5.2. Построение интерполирующих кривых

В дальнейшем будем считать, что множество параметрических узлов задано. Нетрудно видеть, что поставленная задача интерполяции кривой легко сводится к решению двух задач обыкновенной интерполяции: для построения отображения g в (5.48) достаточно построить функции $x : [t_0, t_n] \rightarrow \mathbb{R}$ и $y : [t_0, t_n] \rightarrow \mathbb{R}$, удовлетворяющие условиям

$$\begin{cases} x(t_i) = x_i, \\ y(t_i) = y_i, \end{cases} \quad i = \overline{0, n}. \quad (5.49)$$

Для построения таких x и y можно воспользоваться, например, алгебраической интерполяцией.

Алгебраическая интерполяция

Ввиду того что обе интерполирующие функции x и y строятся по одинаковой сетке узлов, аналог интерполяционного многочлена Лагранжа (5.8) можно записать следующим образом:

$$g(t) = \sum_{i=0}^n \Lambda_i(t) p_i = \sum_{i=0}^n \Lambda_i(t) \begin{bmatrix} x_i \\ y_i \end{bmatrix},$$

где базисная функция Λ_i определяется как обычно:

$$\Lambda_i(t) = \prod_{i \neq j} \frac{t - t_j}{t_i - t_j}.$$

▷₁ Запишите аналоги интерполяционной формулы Ньютона и барицентрических интерполяционных формул.

Сплайновые интерполяционные кривые

Более эффективным инструментом интерполяции кривых являются кубические сплайны. Как обычно, рассмотрим разбиение интервала $[t_0, t_n]$:

$$\Delta = \{[t_{i-1}, t_i]\}_{i=1}^n.$$

Построение сплайновой кривой g сводится к построению двух интерполяционных сплайнов на этом разбиении:

$$x = s^1 \in S_\Delta^3, \quad y = s^2 \in S_\Delta^3$$

по условиям (5.49). Строятся эти сплайны по общей схеме, описанной в подп. 5.4.2. При этом, конечно, нужно задать граничные условия. В итоге мы получим кусочно-кубическую параметрическую функцию вида

$$g(t)|_{t \in \Delta_i} = \begin{bmatrix} \alpha_i^1 \\ \alpha_i^2 \end{bmatrix} + (t - t_i) \begin{bmatrix} \beta_i^1 \\ \beta_i^2 \end{bmatrix} + \frac{1}{2}(t - t_i)^2 \begin{bmatrix} \gamma_i^1 \\ \gamma_i^2 \end{bmatrix} + \frac{1}{6}(t - t_i)^3 \begin{bmatrix} \delta_i^1 \\ \delta_i^2 \end{bmatrix},$$

где $\{\alpha_i^1, \beta_i^1, \gamma_i^1, \delta_i^1\}$ — коэффициенты соответствующих «кусков» кубического сплайна $s^1 = x$, а $\{\alpha_i^2, \beta_i^2, \gamma_i^2, \delta_i^2\}$ — то же самое для $s^2 = y$.

5.5.3. Кривые Безье

Интерактивный дизайн кривой

Рассмотрим теперь задачу *интерактивного дизайна* кривой: требуется построить кривую, обладающую нужной *формой*. Эта задача существенно отличается от рассматриваемой ранее задачи интерполяции кривой тем, что здесь нет строго определенного объекта для приближения. Нужно путем визуального подбора построить кривую, обладающую определенными свойствами (например, имеющую форму кузова автомобиля).

Интерполяция многочленами плохо подходит для решения такой задачи. Это связано с тем, что базисные функции Лагранжа Λ_i хоть и обладают свойством $\Lambda_i(t_j) = \delta_{ij}$, но на всем остальном отрезке интерполяции $[a, b]$ могут принимать значения, намного большие единицы. Это приводит к тому, что небольшое изменение одной точки p_i может сильно изменить значение интерполяционного многочлена в точке t , расположенной достаточно далеко от узла t_i . Этот недостаток можно выразить так: i -я базисная функция Лагранжа Λ_i *не локализована вблизи i -го узла интерполяции*. Гораздо более удобным с этой точки зрения является базис из многочленов Бернштейна.

Базисные многочлены Бернштейна

Определение 5.11. Многочлен

$$B_i^n(t) = \binom{n}{i} t^i (1-t)^{n-i} \quad (5.50)$$

называют i -м (базисным) *многочленом Бернштейна* степени n , $0 \leq i \leq n$. Здесь $\binom{n}{i} = \frac{n!}{i!(n-i)!} = C_n^i$ — биномиальные коэффициенты. При $i < 0$ и $i > n$ полагаем $B_i^n(t) = 0$.

Множество $\{B_i^n\}_{i=0}^n$ образует базис в пространстве многочленов степени не выше n .

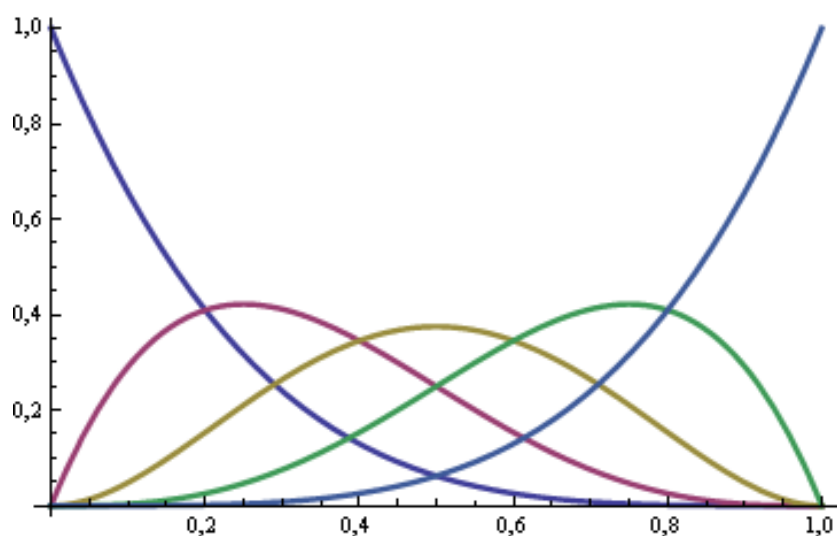


Рис. 5.6. Многочлены Бернштейна степени 4

Графики всех многочленов B_i^4 приведены на рис. 5.6. Многочлены Бернштейна обладают множеством полезных свойств и достаточно широко используются в численном анализе. Наиболее интересны для нас следующие свойства этих многочленов:

- 1) $B_i^n(t) \geq 0 \quad \forall t \in [0, 1], \quad \forall n, i$;
- 2) $B_0^n(0) = B_n^n(1) = 1$;
 $B_i^n(0) = B_i^n(1) = 0 \quad \forall i = 1, n-1$;
- 3) $B_i^n(t) = B_{n-i}^n(1-t)$;
- 4) $\frac{d}{dt} B_i^n(t) = n(B_{i-1}^{n-1}(t) - B_i^{n-1}(t))$;
- 5) B_i^n имеет на отрезке $[0, 1]$ единственный локальный максимум, который достигается в точке i/n (B_i^n локализована в точке i/n);

- 6) $\sum_{i=0}^n B_i^n(t) = 1 \quad \forall t \in \mathbb{R}, \quad n \geq 0;$
 7) $B_i^n(t) = (1-t)B_i^{n-1}(t) + tB_{i-1}^{n-1}(t).$

▷₂ Докажите свойства 4–7.

Кривые Безье

Кривая Безье — это образ единичного отрезка под действием линейной комбинации базисных функций Бернштейна с векторными коэффициентами. Кривые Безье широко применяются в компьютерной графике, с их помощью работают практически все компьютерные векторные шрифты и др.

Определение 5.12. Рассмотрим набор точек плоскости $Q = \{q_i \in \mathbb{R}^2\}_{i=0}^n$, которые в дальнейшем будем называть *контрольными точками*. Кривой Безье называется множество точек

$$\{B(t) \mid t \in [0, 1]\},$$

где

$$B(t) = \sum_{i=0}^n B_i^n(t) q_i. \quad (5.51)$$

Примеры кривых Безье для $n = 1, 2, 3$ изображены на рис. 5.7.

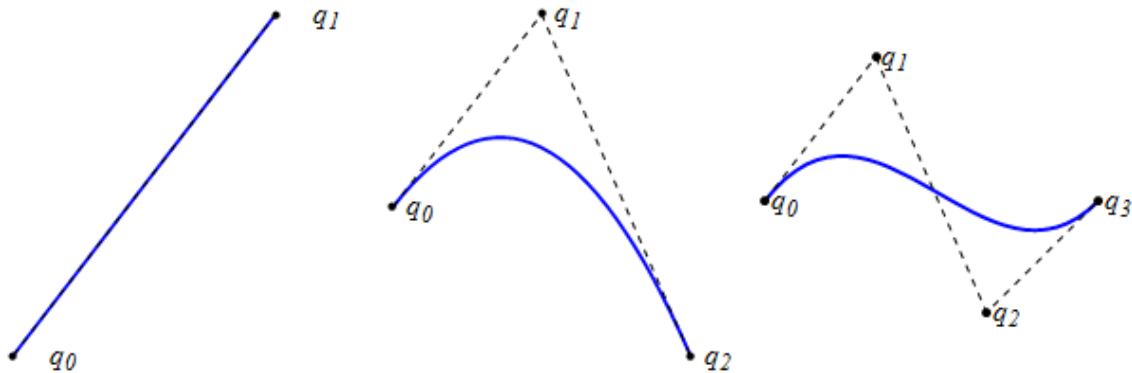


Рис. 5.7. Простейшие кривые Безье

Данный рисунок наглядно иллюстрирует следующие свойства кривых Безье:

- 1) $B(0) = q_0, B(1) = q_n$;
- 2) касательные к кривой в точках $t = 0$ и $t = 1$ коллинеарны векторам $q_1 - q_0$ и $q_n - q_{n-1}$ соответственно;

3) кривая Безье не выходит за пределы выпуклой оболочки множества контрольных точек Q .

Свойство 2 можно сформулировать более конкретно. Согласно свойству 4 многочленов Бернштейна, имеем

$$\begin{aligned}\frac{d}{dt}B(t) &= \sum_{i=0}^n \frac{d}{dt}B_i^n(t)q_i = n \sum_{i=0}^n (B_{i-1}^{n-1}(t) - B_i^{n-1}(t))q_i = \\ &= n \sum_{i=0}^{n-1} B_i^{n-1}(t)(q_{i+1} - q_i).\end{aligned}$$

Отсюда с учетом свойства 2 многочленов Бернштейна получим

$$B'(0) = n(q_1 - q_0), \quad B'(1) = n(q_n - q_{n-1}).$$

5.5.4. Алгоритмы построения кривых Безье

Рассмотрим задачу построения кривой Безье, задаваемой уравнением (5.51), которая сводится к вычислению значений многочлена B в точках $t \in [0, 1]$. Оба рассматриваемых способа основываются на рекуррентных соотношениях для многочленов Бернштейна (свойство 7).

Прямой алгоритм

Вычисление $B(t)$ для данного t можно разбить на два этапа:

- 1) вычисление $b_i = B_i^n(t)$ для всех i от 0 до n ;
- 2) вычисление $B(t) = \sum_{i=0}^n b_i q_i$.

Основная работа выполняется на первом этапе, поэтому рассмотрим его. Непосредственное вычисление b_i по определению (5.50) приводит к вычислительной неустойчивости. Согласно свойству 7, имеем

$$B_i^n(t) = (1-t)B_i^{n-1}(t) + tB_{i-1}^{n-1}(t). \quad (5.52)$$

Для того чтобы вычислить все b_i по такой схеме, необходимо вычислить все $B_j^k(t)$ для $k = \overline{0, n}$, $j = \overline{0, k}$. Осуществляется это последовательным

(слева направо) заполнением следующей треугольной таблицы:

$$\begin{array}{cccccc}
 B_0^0 & B_0^1 & B_0^2 & \dots & B_0^{n-1} & B_0^n \\
 & B_1^1 & B_1^2 & \dots & B_1^{n-1} & B_1^n \\
 & & B_2^2 & \dots & \vdots & \vdots \\
 & & & & B_{n-1}^{n-1} & B_{n-1}^n \\
 & & & & & B_n^n
 \end{array}$$

Хранить всю таблицу не обязательно. Достаточно завести массив b размера $n + 1$ и, начиная с

$$b = (B_0^0, 0, 0, \dots, 0) = (1, 0, 0, \dots, 0),$$

последовательно заполнить все его компоненты, в итоге получив

$$b = (B_0^n, B_1^n, \dots, B_n^n).$$

▷₃ Запишите соответствующий алгоритм.

При построении графика кривой можно в два раза сократить объем вычислений, если учесть симметрию базисных многочленов Бернштейна (свойство 3): один раз вычислив b_i в точке t , можно найти не только $B(t)$, но и

$$B(1 - t) = \sum_{i=0}^n b_{n-i} q_i.$$

Алгоритм de Casteljau

Формулу (5.52) можно применить непосредственно к многочлену (5.51). Если обозначить

$$B(t) = B(q_0, \dots, q_n, t),$$

то из указанных двух формул получим

$$\begin{aligned}
 B(q_0, \dots, q_n, t) &= (1 - t) \sum_{i=0}^n B_i^{n-1} q_i + t \sum_{i=0}^n B_{i-1}^{n-1} q_i = \\
 &= (1 - t) B(q_0, \dots, q_{n-1}, t) + t B(q_1, \dots, q_n, t). \quad (5.53)
 \end{aligned}$$

Таким образом, по рекуррентной схеме (3) мы можем напрямую вычислить $B(t)$, заполняя треугольную таблицу практически аналогично схеме вычисления разделенных разностей. Начнем с массива

$$(B(q_0, t), B(q_1, t), \dots, B(q_n, t)) = (q_0, q_1, \dots, q_n)$$

и постепенно вычислим $B(q_i, \dots, q_{i+k}, t)$, складывая соседние элементы массива с весами $(1 - t)$ и t . Полученная схема вычислений приведена ниже (точку t в обозначении $B(\dots)$ опустим):

$$\begin{array}{ccccccc}
 q_0 & & & & & & \\
 & B(q_0, q_1) & & & & & \\
 q_1 & & B(q_0, q_1, q_2) & & & & \\
 & B(q_1, q_2) & & \dots & & & \\
 q_2 & & B(q_1, q_2, q_3) & & B(q_0, \dots, q_{n-1}) & & \\
 \cdot & B(q_2, q_3) & \cdot & \dots & & B(q_0, \dots, q_n) & \\
 \cdot & \cdot & \cdot & & B(q_1, \dots, q_n) & & \\
 \cdot & \cdot & \cdot & \dots & & & \\
 q_{n-1} & \cdot & B(q_{n-2}, q_{n-1}, q_n) & & & & \\
 & B(q_{n-1}, q_n) & & & & & \\
 q_n & & & & & &
 \end{array}$$

▷₄ Запишите алгоритм.

▷₅ Сравните трудоемкость обоих алгоритмов.

5.6. Среднеквадратичные приближения

Сейчас мы рассмотрим принципиально отличный от интерполяции способ приближения функций — среднеквадратичное приближение. Такие приближения можно строить как по заданному набору точек, так и по заданной функции. Начнем с первого, дискретного, случая.

5.6.1. Дискретный случай

Постановка задачи

Рассмотрим на декартовой плоскости $N + 1$ точек (x_k, y_k) , $k = \overline{0, N}$ (рис. 5.8, слева). Можно считать, что эти точки представляют собой значения некоторой неизвестной функции f , заданные с какой-то погрешностью: $f(x_k) \approx y_k$. Часто в таком случае требуется найти приближение к функции f . Понятно, что интерполяция в данном случае подходит плохо.

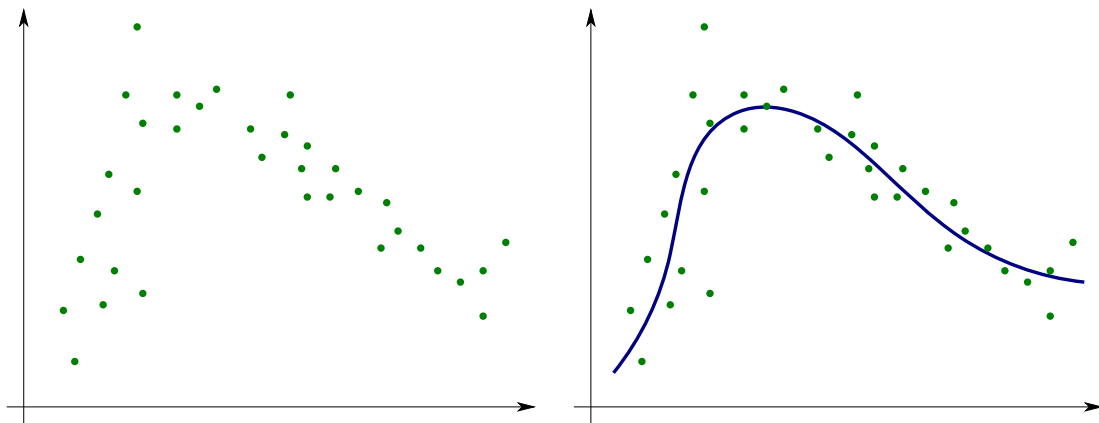


Рис. 5.8. Задача дискретного среднеквадратичного приближения

Вместо нее традиционно применяют метод *среднеквадратичного приближения*, он же *метод наименьших квадратов*. Этот метод состоит в поиске такой функции φ , для которой *сумма квадратов отклонений от y_i в точках x_i* была бы минимальной (рис. 5.8, справа).

Определение 5.13. Задача дискретного среднеквадратичного приближения для данного набора точек (x_k, y_k) , $k = \overline{0, N}$, заключается в построении функции φ вида

$$\varphi(x) = \sum_{i=0}^n \alpha_i \varphi_i(x),$$

для которой выражение

$$\sum_{k=0}^N (y_k - \varphi(x_k))^2 \quad (5.54)$$

принимает минимальное возможное значение. Здесь $\{\varphi_i\}_{i=0}^n$ — некоторый заданный базис функций.

Искомая функция φ однозначно определяется своими коэффициентами $\{\alpha_i\}$, которые представим в виде вектора $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_n)^T$. Запишем выражение (5.54) в виде функционала $\Phi : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$,

$$\Phi(\alpha) = \sum_{k=0}^N \left(y_k - \sum_{i=0}^n \alpha_i \varphi_i(x_k) \right)^2. \quad (5.55)$$

Задача свелась к нахождению вектора α , доставляющего минимум $\Phi(\alpha)$.

Вывод расчетных формул

Рассмотрим внимательнее структуру функционала Φ :

$$\begin{aligned} \Phi(\alpha) &= \sum_{k=0}^N \left(y_k - \sum_{i=0}^n \alpha_i \varphi_i(x_k) \right)^2 = \\ &= \sum_{k=0}^N \left(y_k^2 - 2y_k \sum_{i=0}^n \alpha_i \varphi_i(x_k) + \left(\sum_{i=0}^n \alpha_i \varphi_i(x_k) \right)^2 \right) = \\ &= \sum_{k=0}^N y_k^2 - 2 \sum_{k=0}^N y_k \sum_{i=0}^n \alpha_i \varphi_i(x_k) + \sum_{k=0}^N \sum_{i=0}^n \alpha_i \varphi_i(x_k) \sum_{j=0}^n \alpha_j \varphi_j(x_k) = \\ &= \sum_{k=0}^N y_k^2 - 2 \sum_{i=0}^n \alpha_i \sum_{k=0}^N y_k \varphi_i(x_k) + \sum_{i=0}^n \sum_{j=0}^n \alpha_i \alpha_j \sum_{k=0}^N \varphi_i(x_k) \varphi_j(x_k). \end{aligned}$$

Если ввести обозначения

$$y = (y_0, y_1, \dots, y_N)^T, \quad \beta_i = \sum_{k=0}^N y_k \varphi_i(x_k), \quad \gamma_{ij} = \sum_{k=0}^N \varphi_i(x_k) \varphi_j(x_k),$$

то полученное выражение можно переписать в виде

$$\Phi(\alpha) = (y, y) - 2 \sum_{i=0}^n \beta_i \alpha_i + \sum_{i=0}^n \sum_{j=0}^n \gamma_{ij} \alpha_i \alpha_j. \quad (5.56)$$

Как видим, $\Phi(\alpha)$ представляет собой квадратичный многочлен по каждой переменной α_i , причем $\Phi(\alpha) \geq 0$ при любых $\alpha \in \mathbb{R}^{n+1}$.

▷₁ В каком случае минимальное значение Φ равно нулю?

Значит, единственный глобальный минимум Φ может быть найден как решение системы уравнений

$$\frac{\partial}{\partial \alpha_l} \Phi(\alpha) = 0, \quad l = \overline{0, n}.$$

Непосредственным дифференцированием получим

$$\frac{\partial}{\partial \alpha_l} \Phi(\alpha) = -2\beta_l + 2 \sum_{j=0}^n \gamma_{lj} \alpha_j.$$

▷₂ Выведите эту формулу.

Следовательно, искомый вектор α является решением системы линейных уравнений

$$\sum_{j=0}^n \gamma_{lj} \alpha_j = \beta_l, \quad l = \overline{0, n},$$

которая в матричном виде записывается просто как

$$\Gamma \alpha = \beta, \tag{5.57}$$

$$\begin{aligned} \Gamma &= (\gamma_{ij})_{i,j=0}^n, & \gamma_{ij} &= \sum_{k=0}^N \varphi_i(x_k) \varphi_j(x_k), \\ \beta &= (\beta_0, \beta_1, \dots, \beta_n)^T, & \beta_i &= \sum_{k=0}^N y_k \varphi_i(x_k). \end{aligned} \tag{5.58}$$

Заметим, что матрица Γ представляет собой матрицу скалярных произведений (матрицу Грама) для системы векторов

$$\phi_i = (\varphi_i(x_0), \varphi_i(x_1), \dots, \varphi_i(x_N))^T, \quad i = \overline{0, n}, \tag{5.59}$$

т. е.

$$\gamma_{ij} = (\phi_i, \phi_j),$$

а

$$\beta_i = (y, \phi_i).$$

Таким образом, решение задачи дискретного среднеквадратичного приближения сводится к вычислению вектора неизвестных α , удовлетворяющего СЛАУ (5.57), где матрица Γ и вектор β определяются формулами (5.58). Эта задача однозначно разрешима тогда и только тогда, когда матрица Γ невырождена. В силу вышесказанного для этого необходимо и достаточно, чтобы система векторов $\{\phi_i\}$ (5.59) была линейно независимой. Матрица Γ очевидно является симметричной. Поэтому для решения СЛАУ (5.57) естественно применить метод квадратного корня.

Замечание 5.3. При больших значениях N формулы (5.58) могут давать очень большие по модулю значения коэффициентов матрицы Γ и вектора β , что приведет к большим вычислительным погрешностям. Чтобы этого избежать, следует пронормировать эти коэффициенты, разделив их на N , т. е. фактически решать СЛАУ

$$\frac{1}{N}\Gamma\alpha = \frac{1}{N}\beta.$$

5.6.2. Среднеквадратичное приближение функций

Как уже говорилось, среднеквадратичное приближение можно строить не только для набора точек, но и для аналитически заданной функции f на отрезке $[a, b]$. При этом общая схема решения задачи останется практически без изменений. Главным условием применимости этого метода является *интегрируемость с квадратом* функции f .

Постановка задачи

Определение 5.14. Функция f называется *интегрируемой с квадратом на отрезке $[a, b]$* , если

$$\int_a^b |f(x)|^2 dx < \infty.$$

В этом случае говорят, что f принадлежит пространству $L_2[a, b]$.

Это определение подходит и для комплекснозначных функций f . Мы же в дальнейшем будем считать, что f принимает на $[a, b]$ только вещественные значения.

Определение 5.15. Задача среднеквадратичного приближения интегрируемой с квадратом функции f на отрезке $[a, b]$ заключается в поиске функции φ вида

$$\varphi(x) = \sum_{i=0}^n \alpha_i \varphi_i(x),$$

для которой интеграл

$$\int_a^b (f(x) - \varphi(x))^2 dx \quad (5.60)$$

принимает наименьшее возможное значение. Как обычно, здесь $\{\varphi_i\}$ — заданная система базисных функций.

Замечание 5.4. Иногда вместо (5.60) рассматривают интеграл более общего вида

$$\int_a^b (f(x) - \varphi(x))^2 \rho(x) dx, \quad (5.61)$$

где ρ — некоторая неотрицательная функция, называемая *весовой функцией*.

Вывод расчетных формул

Совершенно аналогично дискретному случаю задача среднеквадратичного приближения функции f сводится к нахождению вектора коэффициентов $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_n)^T$, доставляющего минимум функционала, получаемого подстановкой $\varphi(x) = \sum_{i=0}^n \alpha_i \varphi_i(x)$ в формулу (5.60):

$$\Psi(\alpha) = \int_a^b \left(f(x) - \sum_{i=0}^n \alpha_i \varphi_i(x) \right)^2 dx.$$

Этот функционал имеет ту же структуру, что и Φ для дискретного случая (5.55). Разница лишь в том, что теперь вместо конечной суммы по точкам $\{x_i\}$ мы имеем сумму бесконечную — интеграл. Линейность интеграла позволяет точно так же, как это было сделано выше для $\Phi(\alpha)$, записать $\Psi(\alpha)$ в виде

$$\Psi(\alpha) = \int_a^b f(x)^2 dx - 2 \sum_{i=0}^n \beta_i \alpha_i + \sum_{i=0}^n \sum_{j=0}^n \gamma_{ij} \alpha_i \alpha_j, \quad (5.62)$$

где

$$\beta_i = \int_a^b f(x) \varphi_i(x) dx, \quad \gamma_{ij} = \int_a^b \varphi_i(x) \varphi_j(x) dx, \quad i, j = \overline{0, n}. \quad (5.63)$$

▷₃ Выведите формулы (5.63).

Заметим, что здесь мы использовали те же обозначения (β_i и γ_{ij}), что и ранее, так как они имеют схожий смысл.

Далее по уже известной схеме приравняем $\frac{\partial \Psi(\alpha)}{\partial \alpha_l}$ к нулю и получим следующий результат. *Решение задачи среднеквадратичного приближения функции f сводится к вычислению вектора неизвестных α , удовлетворяющего системе линейных уравнений (5.57),*

$$\Gamma \alpha = \beta,$$

где элементы матрицы $\Gamma = (\gamma_{ij})_{i,j=0}^n$ и вектора $\beta = (\beta_0, \beta_1, \dots, \beta_n)^T$ определяются формулами (5.63). Матрица Γ в данном случае также называется матрицей Грама, так как ее элементы $\gamma_{ij} = \int_a^b \varphi_i(x) \varphi_j(x) dx$ представляют собой скалярные произведения функций φ_i и φ_j . Эта матрица невырождена тогда и только тогда, когда система функций $\{\varphi_i\}$ линейно независима.

▷₄ Получите выражения для β_i и γ_{ij} в случае, когда функционал Ψ определяется формулой (5.61).

5.6.3. Базисные функции

Обсудим теперь вопрос выбора базисных функций $\{\varphi_i\}$ для построения среднеквадратичного приближения.

Степенной базис

Если мы хотим, чтобы приближающая функция φ была алгебраическим многочленом, самый простой базис в этом случае имеет вид

$$\varphi_i(x) = x^i, \quad i = \overline{0, n}. \quad (5.64)$$

Матрица Грама Γ для этого базиса вычисляется легко (предположим, что $[a, b] = [0, 1]$):

$$\gamma_{ij} = \int_0^1 x^{i+j} dx = \frac{1}{i+j+1}, \quad \text{т. е.}$$

$$\Gamma = H_{n+1} = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \cdots & \frac{1}{n+1} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \cdots & \frac{1}{n+2} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \cdots & \frac{1}{n+3} \\ \dots & \dots & \dots & \dots & \dots \\ \frac{1}{n+1} & \frac{1}{n+2} & \frac{1}{n+3} & \cdots & \frac{1}{2n+1} \end{bmatrix}.$$

Это знаменитая матрица Гильберта, классический пример плохо обусловленной матрицы. Ее число обусловленности очень быстро растет с ростом n , поэтому на практике базис (5.64) применяется лишь при небольших n .

▷₅ Вычислите матрицу Грама для степенного базиса в случае произвольного отрезка $[a, b]$.

▷₆ Запишите вид соответствующего вектора β .

Вместо базиса (5.64) для более эффективного вычисления среднеквадратичного приближения полиномами используются различные системы *ортogonalных многочленов*, с которыми мы познакомимся в п. 6.4.

Классический тригонометрический базис

На практике часто используется тригонометрический базис на отрезке $[-\pi, \pi]$:

$$\{\varphi_i(x)\}_{i=0}^n = \{1, \cos x, \sin x, \cos 2x, \sin 2x, \dots, \cos Kx, \sin Kx\},$$

где $n = 2K$. Данный базис замечателен тем, что он является *ортogonalным*, т. е. базисом, для которого матрица Грама диагональна (в случае единичной весовой функции ρ). Действительно, при таком выборе φ_i имеем

$$\gamma_{ij} = (\varphi_i, \varphi_j) = \int_{-1}^1 \varphi_i(x) \varphi_j(x) dx = \begin{cases} 0, & i \neq j, \\ 2\pi, & i = j = 0, \\ \pi, & i = j \neq 0. \end{cases}$$

Следовательно, решение системы $\Gamma\alpha = \beta$ для нахождения коэффициентов среднеквадратичного приближения можно выписать сразу:

$$\begin{aligned}\alpha_0 &= \frac{1}{2\pi}\beta_0 = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x)dx, \\ \alpha_{2k-1} &= \frac{1}{\pi}\beta_{2k-1} = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos kx dx, \\ \alpha_{2k} &= \frac{1}{\pi}\beta_{2k} = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin kx dx, \quad k = \overline{1, K}.\end{aligned}$$

Полученное таким образом приближение φ является отрезком *ряда Фурье* для функции f . Ряды Фурье широко используются в прикладной математике, подробности можно найти, например, в [9, гл. IV].

Фундаментальный кусочно-линейный базис

Рассмотрим разбиение отрезка $[a, b] = [x_0, x_n]$ на n частей:

$$\Delta = \{\Delta_i\}_{i=1}^n, \quad \Delta_i = [x_{i-1}, x_i]$$

и S_{Δ}^1 — множество сплайнов первого порядка на этом разбиении. Напомним, что элементами этого множества являются непрерывные кусочно-линейные функции. Предположим, что необходимо найти среднеквадратичное приближение к функции f среди элементов S_{Δ}^1 .

Для решения этой задачи необходимо задать базис пространства S_{Δ}^1 . Наиболее популярным является фундаментальный базис, i -я функция в котором принимает значение 1 в узле x_i и 0 во всех остальных узлах. Этого описания достаточно для определения базиса, так как любая функция из множества S_{Δ}^1 однозначно определяется своими значениями в узлах $\{x_i\}$. Аналитические формулы для базисных функций (обозначим их $\{\psi_i\}$) выглядят следующим образом:

$$\begin{aligned}\psi_0(x) &= \begin{cases} \frac{x_1 - x}{h_1}, & x \in \Delta_1, \\ 0, & x \notin \Delta_1, \end{cases} & \psi_n(x) &= \begin{cases} \frac{x - x_{n-1}}{h_n}, & x \in \Delta_n, \\ 0, & x \notin \Delta_n, \end{cases} \\ \psi_i(x) &= \begin{cases} \frac{x - x_{i-1}}{h_i}, & x \in \Delta_i, \\ \frac{x_{i+1} - x}{h_{i+1}}, & x \in \Delta_{i+1}, \\ 0, & x \notin \Delta_i \cup \Delta_{i+1}, \end{cases} & i &= \overline{1, n-1}.\end{aligned}$$

Напомним обозначения: здесь $\Delta_i = [x_{i-1}, x_i]$, $h_i = x_i - x_{i-1}$.

Графики $\{\psi_i\}$ для частного случая $n = 3$ приведены на рис. 5.9.

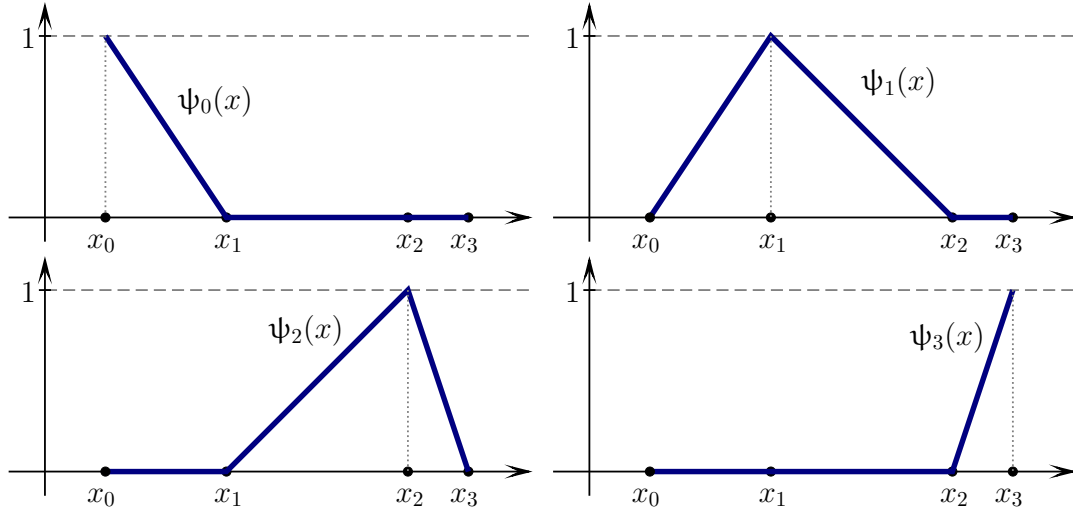


Рис. 5.9. Базис пространства кусочно-линейных функций на сетке

Интересно вычислить матрицу Грама для построенного базиса. Эта матрица будет трехдиагональной, так как при $|i - j| \geq 2$ имеем $\psi_i(x)\psi_j(x) = 0$ для всех x . Прямым вычислением получим

$$\Gamma = \begin{bmatrix} d_0 & c_1 & & & \\ c_1 & d_1 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-1} & d_{n-1} & c_n \\ & & & c_n & d_n \end{bmatrix}, \quad (5.65)$$

где

$$\begin{aligned} d_i &= \frac{h_i + h_{i+1}}{3}, \quad i = \overline{1, n-1}, \quad d_0 = \frac{h_1}{3}, \quad d_n = \frac{h_n}{3}, \\ c_i &= \frac{h_i}{6}, \quad i = \overline{1, n}. \end{aligned} \quad (5.66)$$

▷7 Выведите эти формулы.

5.7. Приближение поверхностей

5.7.1. Интерполяция на прямоугольнике

Задача приближения функции одной переменной, которую мы рассматривали до сих пор, естественным образом обобщается на случай функций нескольких переменных. Основное внимание в дальнейшем мы уделим задаче интерполяции функций двух переменных на прямоугольной области.

Определение 5.16. Рассмотрим на декартовой плоскости прямоугольник $\Pi = [a, b] \times [c, d]$. На этом прямоугольнике зададим множество точек — *прямоугольную сетку*

$$(x_i, y_j), \quad i = \overline{0, n}, \quad j = \overline{0, m},$$

порождаемую одномерными сетками на $[a, b]$ и $[c, d]$ соответственно (рис. 5.10):

$$a = x_0 < x_1 < \dots < x_n = b, \quad \text{и} \quad c = y_0 < y_1 < \dots < y_m = d.$$

Задача интерполяции функции $f : \Pi \rightarrow \mathbb{R}$ на данной сетке узлов заключается в построении функции $\varphi : \Pi \rightarrow \mathbb{R}$, такой, что

$$\varphi(x_i, y_j) = z_{ij} = f(x_i, y_j) \quad i = \overline{0, n}, \quad j = \overline{0, m}.$$

Вид функции φ , как правило, задается заранее.

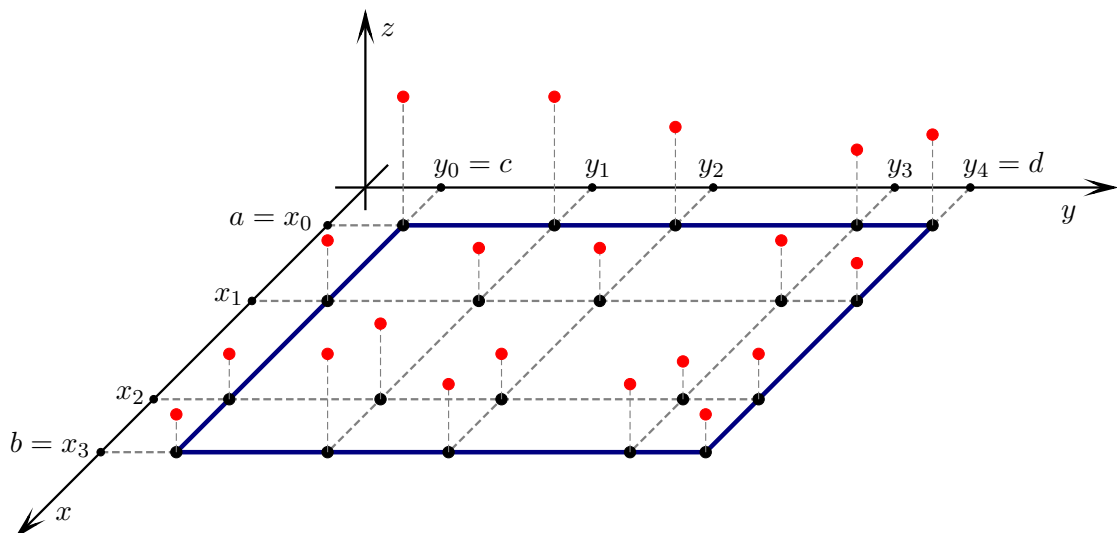


Рис. 5.10. Двумерная интерполяция на прямоугольной сетке

Преимущество такой постановки задачи в ее простоте: для интерполяции по прямоугольнику легко адаптируются одномерные методы интерполяции, которые мы успели рассмотреть.

5.7.2. Полиномиальная интерполяция

Рассмотрим случай, когда искомая интерполирующая функция φ является алгебраическим многочленом двух переменных, а точнее многочленом степени n по переменной x и степени m по y :

$$\varphi(x, y) = P_{n,m}(x, y) = \sum_{i=0}^n \sum_{j=0}^m \alpha_{ij} x^i y^j.$$

Такую задачу интерполяции можно решить «методом грубой силы»: составить систему из $(n+1)(m+1)$ линейных уравнений для нахождения неизвестных коэффициентов $\{\alpha_{ij}\}$. Уравнения будут иметь вид

$$\sum_{i=0}^n \sum_{j=0}^m \alpha_{ij} x_p^i y_q^j = z_{pq}, \quad p = \overline{0, n}, \quad q = \overline{0, m}.$$

Но это слишком нерационально. Гораздо проще обобщить формулу интерполяционного многочлена Лагранжа (5.8) на случай двух переменных.

Сначала построим одномерные базисные функции Лагранжа (5.7) отдельно для узлов $\{x_i\}_{i=0}^n$ и $\{y_j\}_{j=0}^m$:

$$\begin{aligned} \Lambda_i^1(x) &= \prod_{p \neq i} \frac{x - x_p}{x_i - x_p}, \quad i = \overline{0, n}, \\ \Lambda_j^2(y) &= \prod_{q \neq j} \frac{y - y_q}{y_j - y_q}, \quad j = \overline{0, m}. \end{aligned}$$

Теперь рассмотрим базисные многочлены двух переменных

$$\Lambda_{ij}(x, y) = \Lambda_i^1(x) \Lambda_j^2(y).$$

Это *фундаментальный базис двумерной алгебраической интерполяции*, так как по построению имеем

$$\Lambda_{ij}(x_p, y_q) = \begin{cases} 1, & \text{если } p = i \text{ и } q = j \\ 0, & \text{в противном случае.} \end{cases}$$

Следовательно, многочлен

$$P_{n,m}(x, y) = \sum_{i=0}^n \sum_{j=0}^m z_{ij} \Lambda_{ij}(x, y) \quad (5.67)$$

удовлетворяет условиям интерполяции

$$P_{n,m}(x_i, y_j) = z_{ij}, \quad \forall i = \overline{0, n}, \quad j = \overline{0, m},$$

т. е. является искомым интерполяционным многочленом в форме Лагранжа.

▷₁ Постройте аналог барицентрической формулы Лагранжа для двумерного случая.

5.7.3. Интерполяция билинейными сплайнами

Теперь рассмотрим обобщение понятия линейного сплайна на случай двух переменных. Начнем с обозначений для отрезков:

$$\begin{aligned} \Delta_i^1 &= [x_{i-1}, x_i], \quad i = \overline{1, n}, \\ \Delta_j^2 &= [y_{j-1}, y_j], \quad j = \overline{1, m}. \end{aligned}$$

Соответствующие разбиения обозначим $\Delta^1 = \{\Delta_i^1\}_{i=1}^n$, $\Delta^2 = \{\Delta_j^2\}_{j=1}^m$. Эти разбиения порождают разбиение прямоугольника $\Pi = [a, b] \times [c, d]$ на множество прямоугольников

$$\Delta = \Delta^1 \times \Delta^2 = \{\Delta_{ij} = \Delta_i^1 \times \Delta_j^2\}_{i=1, j=1}^{n, m}.$$

Определение 5.17. *Билинейным сплайном* на разбиении $\Delta = \Delta^1 \times \Delta^2$ называется непрерывная функция $s : \Pi \rightarrow \mathbb{R}$, которая на каждом прямоугольнике Δ_{ij} является *билинейной функцией* (рис. 5.11):

$$s(x, y) \Big|_{(x,y) \in \Delta_{ij}} = s_{ij}(x, y) = \alpha_{ij}^{00} + \alpha_{ij}^{10}x + \alpha_{ij}^{01}y + \alpha_{ij}^{11}xy. \quad (5.68)$$

Заметим, что можно дать и другое, эквивалентное, определение билинейного сплайна: это функция двух переменных s такая, что при любом фиксированном $x \in [a, b]$ выражение $s(x, y)$ как функция от y является сплайном первого порядка на разбиении Δ^2 . А при любом фиксированном $y \in [c, d]$ выражение $s(x, y)$ как функция от x является сплайном первого порядка на разбиении Δ^1 .

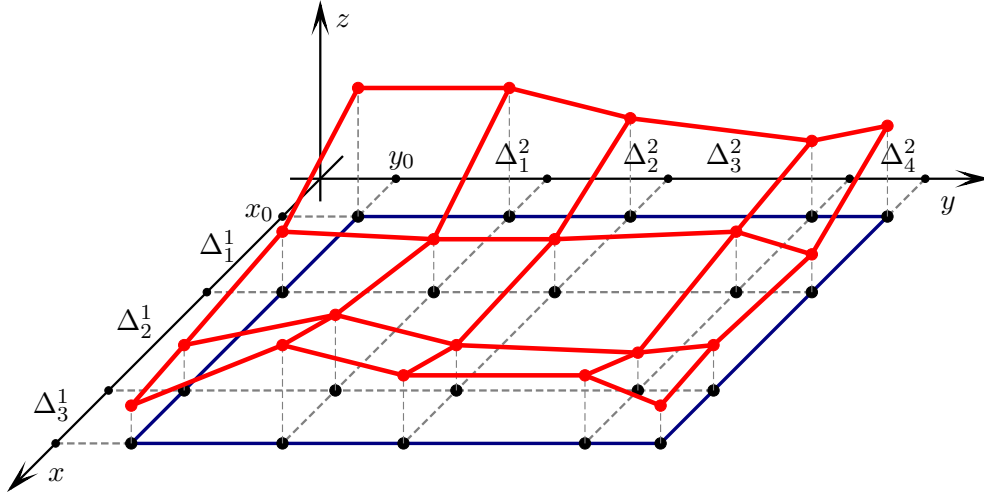


Рис. 5.11. Билинейный сплайн

Построение интерполяционного билинейного сплайна, удовлетворяющего условиям

$$s(x_i, y_j) = z_{ij},$$

не представляет проблем. На каждом прямоугольнике Δ_{ij} «кусоч» сплайна s_{ij} (5.68) можно построить, используя формулу Лагранжа (5.67) по очевидным условиям

$$s_{ij}(x_{i-1}, y_{j-1}) = z_{i-1,j-1},$$

$$s_{ij}(x_i, y_{j-1}) = z_{i,j-1},$$

$$s_{ij}(x_{i-1}, y_j) = z_{i-1,j},$$

$$s_{ij}(x_i, y_j) = z_{i,j}.$$

▷₂ Выведите формулы для вычисления $s_{ij}(x, y)$.

5.7.4. Интерполяция бикубическими сплайнами

Немного сложнее обстоят дела с кубическими сплайнами двух переменных — бикубическими сплайнами, которые будут сейчас рассмотрены. Для начала дадим определение с учетом введенных выше обозначений.

Определение 5.18. Бикубическим сплайном, определенном на разбиении $\Delta = \Delta^1 \times \Delta^2$, называется функция $s : \Pi \rightarrow R$, удовлетворяющая следующим двум условиям:

- при любом фиксированном значении $x \in [a, b]$ выражение $s(x, y)$ как функция переменной y является кубическим сплайном на разбиении Δ^2 ;

- при любом фиксированном значении $y \in [c, d]$ выражение $s(x, y)$ как функция переменной x является кубическим сплайном на разбиении Δ^1 .

Из этого определения следует, что:

- 1) на каждом прямоугольнике Δ_{ij} бикубический сплайн s задается многочленом третьей степени по каждой переменной;
- 2) первая и вторая производные бикубического сплайна, включая смешанные, непрерывны на всем прямоугольнике Π .

Эти два свойства можно считать эквивалентным определением бикубического сплайна (рис. 5.12).

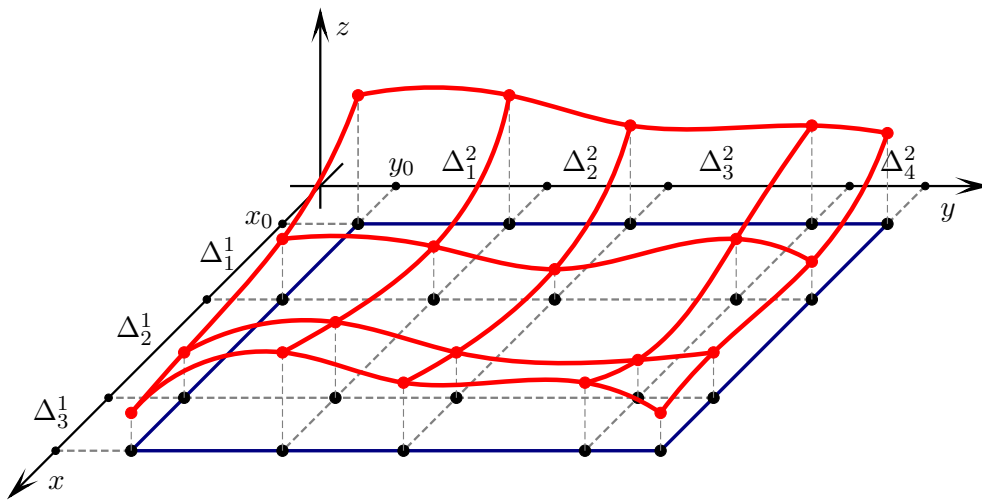


Рис. 5.12. Бикубический сплайн

Приступим теперь к построению интерполяционного бикубического сплайна s , удовлетворяющего условиям

$$s(x_i, y_j) = z_{ij}.$$

Для решения этой задачи достаточно уметь строить одномерные кубические сплайны (см. подп. 5.4.2).

Зафиксируем произвольное $y \in [c, d]$. Тогда по определению $s(x, y)$ представляет собой обычный кубический сплайн по переменной x , который можно записать, как мы делали это ранее:

$$s(x, y) \Big|_{x \in \Delta_i^1} = \alpha_i(y) + \beta_i(y)(x - x_i) + \frac{\gamma_i(y)}{2}(x - x_i)^2 + \frac{\delta_i(y)}{6}(x - x_i)^3. \quad (5.69)$$

Здесь коэффициенты $\alpha_i, \beta_i, \gamma_i, \delta_i, i = \overline{1, n}$, очевидно зависят от зафиксированной нами точки y . Более того, из определения бикубического сплайна следует, что эти коэффициенты, как функции переменной y , являются кубическими сплайнами на разбиении Δ^2 :

$$\alpha_i, \beta_i, \gamma_i, \delta_i \in S_{\Delta^2}^3, \quad i = \overline{1, n}.$$

Как только мы определим эти сплайны, а их в общей сложности $4n$ штук, искомый бикубический сплайн будет построен.

Необходимые для этого условия получаются из (5.69) подстановкой $y = y_j, j = \overline{0, m}$. В обозначениях

$$\alpha_i(y_j) = \alpha_{ij}, \quad \beta_i(y_j) = \beta_{ij}, \quad \delta_i(y_j) = \delta_{ij}, \quad \gamma_i(y_j) = \gamma_{ij} \quad (5.70)$$

эти условия для каждого j имеют вид

$$\begin{aligned} s(x, y_j) \Big|_{x \in \Delta_i^1} &= u_j(x) \Big|_{x \in \Delta_i^1} = \\ &= \alpha_{ij} + \beta_{ij}(x - x_i) + \frac{\gamma_{ij}}{2}(x - x_i)^2 + \frac{\delta_{ij}}{6}(x - x_i)^3, \quad i = \overline{1, n}. \end{aligned}$$

Здесь u_j — кубический сплайн на разбиении Δ^1 . Его коэффициенты $\{\alpha_{ij}, \beta_{ij}, \gamma_{ij}, \delta_{ij}\}_{i=1}^n$ легко найти, так как по условию u_j должен удовлетворять условиям интерполяции

$$u_j(x_i) = z_{ij}, \quad i = \overline{0, n}.$$

При этом, как и в одномерном случае, для однозначного определения сплайна нужно будет задать дополнительные граничные условия. Это могут быть, в принципе, любые условия из рассмотренных нами ранее для кубических сплайнов. Вычисляя сплайны u_j для $j = \overline{0, m}$ по схеме, описанной в подп. 5.4.2, мы находим $\{\alpha_{ij}, \beta_{ij}, \gamma_{ij}, \delta_{ij}\}_{i=1}^n$.

После этого остается вычислить искомые сплайны $\alpha_i, \beta_i, \gamma_i, \delta_i$ по условиям интерполяции (5.70) для j от 0 до m , а также по соответствующим граничным условиям.

▷3 Укажите количество и размерность СЛАУ, которые нужно решить для построения бикубического сплайна таким способом.

Замечание 5.5. Строить бикубический сплайн можно и в другой последовательности, если вместо y в самом начале зафиксировать переменную x .

5.7.5. Поверхности Безье

Теперь от интерполяции перейдем к задаче интерактивного дизайна поверхностей, используя обобщение кривой Безье — поверхность Безье.

Определение 5.19. Рассмотрим набор точек в трехмерном пространстве $Q = \{q_{ij} \in \mathbb{R}^3\}_{i=0,j=0}^{n,m}$, которые в дальнейшем будем называть *контрольными точками*. Поверхностью Безье называется множество точек

$$\{B(u, v) \mid u, v \in [0, 1]\},$$

где

$$B(u, v) = \sum_{i=0}^n \sum_{j=0}^m B_i^n(u) B_j^m(v) q_{ij}. \quad (5.71)$$

Здесь B_i^n и B_j^m — многочлены Бернштейна (5.50).

Контрольные точки $\{q_{ij}\}$ играют здесь ту же роль, что и в случае кривой Безье: они контролируют форму поверхности. В частности, поверхность Безье проходит через угловые точки q_{00} , q_{n0} , q_{m0} и q_{mn} , а также не выходит за пределы выпуклой оболочки контрольных точек.

Для построения поверхностей Безье используются те же вычислительные алгоритмы, что и для кривых.

▷₄ Запишите адаптацию алгоритма de Casteljau для поверхностей Безье.

Глава 6

Приближенное вычисление интегралов

- 6.1. Интерполяционные квадратурные формулы
- 6.2. Остаток квадратурных формул
- 6.3. Квадратурные формулы Ньютона – Котеса
- 6.4. Квадратурные формулы Гаусса
- 6.5. Составные квадратурные формулы

6.1. Интерполяционные квадратурные формулы

6.1.1. Постановка задачи численного интегрирования

Рассмотрим задачу вычисления интеграла

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

где f — интегрируемая по Риману функция на отрезке $[a, b]$. Как известно, далеко не всегда значение интеграла может быть выражено через элементарные функции, поэтому необходимы методы приближенного интегрирования, которые будут рассматриваться далее.

Теоретически для приближенного вычисления $\int_a^b f(x)dx$ можно воспользоваться определением интеграла Римана как предела интегральных сумм вида

$$\sum_i f(\xi_i)\Delta x_i,$$

где $\{x_i\}$ — разбиение отрезка интегрирования; $\Delta x_i = x_i - x_{i-1}$; $\xi_i \in [x_{i-1}, x_i]$. Но такой способ, как правило, сходится очень медленно. Гораздо эффективнее *приблизить подынтегральную функцию некоторой функцией φ , интеграл от которой вычисляется легко, и заменить интеграл от f на интеграл от φ* . Формулы для приближенного вычисления интегралов называются *квадратурными формулами (КФ)*.

6.1.2. Простейшие квадратурные формулы

Формулы прямоугольников

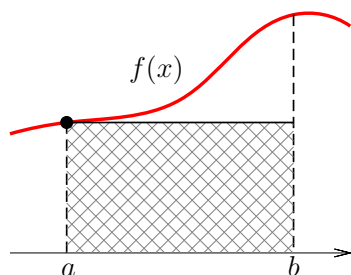
Самые простые квадратурные формулы получаются при замене подынтегральной функции f на константу, а точнее на значение $f(x_0)$, где x_0 — некоторая точка из интервала интегрирования $[a, b]$. Площадь криволинейной трапеции, которой по определению равен искомый интеграл $\int_a^b f(x)dx$, приближенно равна площади прямоугольника $[a, b] \times [0, f(x_0)]$,

т. е.

$$\int_a^b f(x)dx \approx f(x_0)(b-a).$$

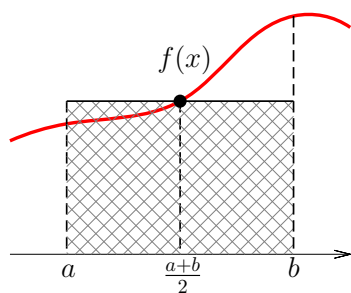
Наиболее известны три квадратурные формулы такого типа:

- *формула левых прямоугольников:*



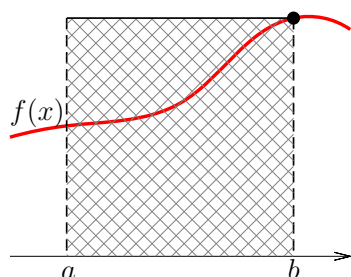
$$\int_a^b f(x)dx \approx (b-a)f(a); \quad (6.1)$$

- *формула средних прямоугольников:*



$$\int_a^b f(x)dx \approx (b-a)f\left(\frac{a+b}{2}\right); \quad (6.2)$$

- *формула правых прямоугольников:*

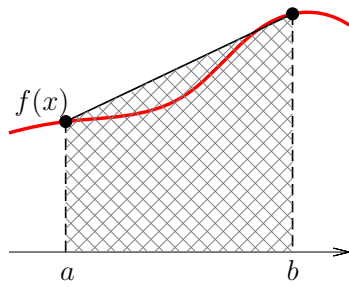


$$\int_a^b f(x)dx \approx (b-a)f(b). \quad (6.3)$$

Квадратурная формула трапеций

Если график подынтегральной функции приблизить секущей, проходящей через точки $f(a)$ и $f(b)$, что при интегрировании эквивалентно замене площади криволинейной трапеции на площадь трапеции обычно-

венной, то в результате получится *квадратурная формула трапеций*:



$$\int_a^b f(x)dx \approx (b-a) \frac{f(a) + f(b)}{2}. \quad (6.4)$$

Эта формула есть не что иное, как школьная формула для вычисления площади трапеции.

6.1.3. Определение

Перейдем теперь от частных случаев к общему определению квадратурной формулы. Начнем с того, что запишем интеграл в более общем виде

$$S(f) = \int_a^b f(x)\rho(x)dx. \quad (6.5)$$

Здесь ρ — некоторая неотрицательная весовая функция, которая может принимать нулевые значения лишь в конечном числе точек. Отображение $S : f \mapsto \int_a^b f(x)\rho(x)dx$ представляет собой линейный функционал.

Определение 6.1. На отрезке $[a, b]$ рассмотрим множество из $n + 1$ попарно различных точек x_0, x_1, \dots, x_n . Формула для приближенного вычисления интеграла $S(f)$ вида

$$\int_a^b f(x)\rho(x)dx \approx \sum_{i=0}^n A_i f(x_i) \quad (6.6)$$

называется *квадратурной формулой*. Вещественные числа $\{A_i\}_{i=0}^n$ называются *коэффициентами квадратурной формулы*, а точки $\{x_i\}_{i=0}^n$ — ее *узлами*. Функционал, стоящий в правой части квадратурной формулы,

$$Q_n(f) = \sum_{i=0}^n A_i f(x_i),$$

называется *квадратурной суммой*.

Таким образом, любая квадратурная формула может быть кратко записана в виде

$$S(f) \approx Q_n(f).$$

6.1.4. Интерполяционные квадратурные формулы

Как уже говорилось ранее, традиционный способ построения квадратурных формул заключается в замене подынтегральной функции f на легко интегрируемую функцию φ , т. е.

$$Q_n(f) = S(\varphi).$$

Если φ является интерполяционным многочленом, то получаемая в результате квадратурная формула называется *интерполяционной*.

Определение 6.2. Квадратурная формула вида

$$S(f) \approx S(P_n),$$

где P_n — интерполяционный многочлен для f , построенный по узлам $\{x_i\}_{i=0}^n$, называется *интерполяционной квадратурной формулой*.

Запишем интерполяционную квадратурную формулу в каноническом виде (6.6), используя представление интерполяционного многочлена в форме Лагранжа:

$$P_n(x) = \sum_{i=0}^n \Lambda_i(x) f(x_i).$$

Здесь Λ_i — базисные функции (5.7). Тогда

$$S(P_n) = S\left(\sum_{i=0}^n \Lambda_i f(x_i)\right) = \sum_{i=0}^n S(\Lambda_i) f(x_i) = \sum_{i=0}^n A_i f(x_i) = Q_n(f).$$

Таким образом, коэффициенты *интерполяционной квадратурной формулы с узлами $\{x_i\}_{i=0}^n$ могут быть вычислены по правилу*

$$\begin{aligned} A_i &= S(\Lambda_i) = \int_a^b \Lambda_i(x) \rho(x) dx, \\ \Lambda_i(x) &= \prod_{j \neq i} \frac{x - x_j}{x_i - x_j}. \end{aligned} \tag{6.7}$$

Очевидно, что уже знакомые формулы прямоугольников и трапеций являются интерполяционными и соответствуют весовой функции $\rho(x) \equiv 1$.

▷₁ Получите формулу трапеций, используя правило (6.7).

6.2. Остаток квадратурных формул

Определение 6.3. Остатком квадратурной формулы $S(f) \approx Q_n(f)$ называется функционал

$$R_n(f) = S(f) - Q_n(f), \quad (6.8)$$

или, в развернутой форме,

$$R_n(f) = \int_a^b f(x)\rho(x)dx - \sum_{i=0}^n A_i f(x_i).$$

Величина остатка — естественная характеристика точности квадратурной формулы.

6.2.1. Остаток интерполяционных квадратурных формул

Имея выражение для остатка полиномиального интерполирования

$$r_n = f - P_n,$$

где P_n — интерполяционный многочлен для f , легко получить выражение для остатка интерполяционных квадратур. Действительно:

$$R_n(f) = S(f) - Q_n(f) = S(f) - S(P_n) = S(f - P_n) = S(r_n),$$

т. е. *остаток интерполяционной квадратурной формулы равен интегралу от остатка интерполяции*. Отсюда, используя формулу (5.28), которая имеет вид

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

для достаточно гладких f получим

$$R_n(f) = \frac{1}{(n+1)!} \int_a^b f^{(n+1)}(\xi) \omega_{n+1}(x) \rho(x) dx. \quad (6.9)$$

Напомним, что здесь $\omega_{n+1}(x) = (x - x_0) \dots (x - x_n)$, а величина ξ зависит от x , поэтому выносить $f^{(n+1)}(\xi)$ за знак интеграла, вообще говоря,

нельзя. Однако в некоторых случаях это возможно. Например, если многочлен ω_{n+1} сохраняет знак на $[a, b]$, то применима одна из теорем о среднем для интеграла Римана.

Теорема 6.1 (о среднем). *Если функция f непрерывна на отрезке $[a, b]$, а функция g интегрируема и знакопостоянна на $[a, b]$, то существует такая точка $\eta \in [a, b]$, что*

$$\int_a^b f(x)g(x)dx = f(\eta) \int_a^b g(x)dx.$$

Таким образом, если многочлен ω_{n+1} знакопостоянен, то по теореме о среднем из (6.9) получим

$$R_n(f) = \frac{f^{(n+1)}(\eta)}{(n+1)!} \int_a^b \omega_{n+1}(x)\rho(x)dx, \quad \eta \in [a, b]. \quad (6.10)$$

Здесь также использован тот факт, что весовая функция ρ знакопостоянна по условию.

6.2.2. Остаток простейших квадратурных формул

Применим полученные формулы для исследования остатка рассмотренных в подп. 6.1.2 простейших квадратурных формул. Все они являются интерполяционными с единичной весовой функцией. Начнем с формул прямоугольников.

Остаток формул левых и правых прямоугольников

Формула (6.1) — это интерполяционная квадратурная формула с $n = 0$ (подынтегральная функция приближается многочленом нулевой степени) и $x_0 = a$:

$$S(f) \approx Q_0(f) = (b-a)f(a).$$

В силу знакопостоянства многочлена $\omega_1(x) = x - a$ остаток этой квадратурной формулы может быть представлен в виде (6.10). В результате получается

$$R_0(f) = \frac{(b-a)^2}{2} f'(\eta), \quad \eta \in [a, b]. \quad (6.11)$$

▷₁ Запишите остаток квадратурной формулы правых прямоугольников (6.3).

Остаток формулы средних прямоугольников

Формула (6.2) наиболее интересна. Перепишем ее в виде

$$S(f) \approx Q_0(f) = (b-a)f\left(\frac{a+b}{2}\right).$$

В этом случае многочлен $\omega_1 = x - (a+b)/2$ уже не сохраняет знак на отрезке интегрирования, т. е. формула (6.10) не работает. Однако хорошее представление для остатка можно получить следующим образом. Предположим, что f'' существует и непрерывна. Тогда по формуле Тейлора имеем

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(\xi)}{2}(x - x_0)^2, \quad \xi \in [a, b].$$

Вычислим $R_0(f) = S(f - P_0)$, где $P_0(x) \equiv f(x_0)$, $x_0 = (a+b)/2$:

$$\begin{aligned} S(f - P_0) &= \int_a^b \left(f'(x_0)(x - x_0) + \frac{f''(\xi)}{2}(x - x_0)^2 \right) dx = \\ &= \int_a^b \frac{f''(\xi)}{2}(x - x_0)^2 dx = \frac{f''(\eta)}{2} \int_a^b (x - x_0)^2 dx, \end{aligned}$$

откуда

$$R_0(f) = \frac{(b-a)^3}{24} f''(\eta). \quad (6.12)$$

▷₂ Какой вид будет иметь остаток формулы средних прямоугольников в случае, когда f'' не существует, т. е. когда работает только формула (6.9)?

Остаток формулы трапеций

Для формулы трапеций (6.4), которую запишем в виде

$$S(f) \approx Q_1(f) = \frac{b-a}{2} f(a) + \frac{b-a}{2} f(b),$$

имеем $n = 1$, $x_0 = a$ и $x_1 = b$. Многочлен $\omega_2(x) = (x-a)(x-b)$ знакопостоянен на $[a, b]$, поэтому снова имеем право применить формулу (6.10) для остатка. В результате получим

$$R_1(f) = -\frac{(b-a)^3}{12} f''(\eta), \quad \eta \in [a, b]. \quad (6.13)$$

▷₃ Выведите эту формулу. При взятии интеграла будет удобно сделать замену переменных $x = a + (b - a)t$.

Если сравнить формулы (6.12) и (6.13), видно, что формула средних прямоугольников, которая требует лишь одного вычисления функции f , точнее формулы трапеций при условии достаточной гладкости f .

6.3. Квадратурные формулы Ньютона – Котеса

6.3.1. Алгебраическая степень точности

Помимо величины остатка существует еще одна характеристика квадратурных формул — алгебраическая степень точности.

Определение 6.4. Рассмотрим квадратурную формулу $S(f) \approx Q_n(f)$. Если для любого многочлена p степени m и ниже выполняется тождество

$$S(p) = Q_n(p),$$

и при этом существует такой многочлен \tilde{p} степени $m + 1$, что

$$S(\tilde{p}) \neq Q_n(\tilde{p}),$$

то говорят, что алгебраическая степень точности (АСТ) такой формулы равна m .

Замечание 6.1. В силу линейности функционалов S и Q_n можно дать эквивалентное определение АСТ. Любой многочлен степени m является линейной комбинацией одночленов

$$\chi_j(x) = x^j, \quad j = \overline{0, m},$$

поэтому квадратурная формула имеет АСТ m , если и только если она точно интегрирует все многочлены χ_j :

$$S(\chi_j) = Q_n(\chi_j) \quad \forall j = \overline{0, m},$$

и при этом

$$S(\chi_{m+1}) \neq Q_n(\chi_{m+1}).$$

▷₁ Докажите это.

Поскольку интерполяционные квадратурные формулы получаются заменой подынтегральной функции на интерполяционный многочлен степени n , а все многочлены степени $m \leq n$ при этом будут интерполироваться точно, получается следующий простой, но важный результат.

Теорема 6.2. Квадратурная формула с $(n + 1)$ узлом является интерполяционной, если и только если она имеет алгебраическую степень точности не менее n .

Доказательство. Если КФ является интерполяционной, то $Q_n(f) = S(P_n)$, где P_n — интерполяционный многочлен степени n для f . Если f является многочленом степени $m \leq n$, то имеем $P_n = f$ и $Q_n(f) = S(f)$ — формула точна для всех многочленов степени n и ниже.

Пусть квадратурная формула

$$S(f) \approx \sum_i A_i f(x_i)$$

имеет АСТ не менее n . Покажем, что ее коэффициенты находятся по формуле (6.7): $A_i = S(\Lambda_i)$, где Λ_i — фундаментальная базисная функция Лагранжа (5.7). Это и будет означать, что формула интерполяционная. Нет ничего проще: так как $\deg \Lambda_i = n$, имеем

$$S(\Lambda_i) = Q_n(\Lambda_i) = \sum_{j=0}^n A_j \Lambda_i(x_j) = A_i. \quad \square$$

▷₂ Найдите АСТ всех трех формул прямоугольников и формулы трапеций.

6.3.2. Симметричные квадратурные формулы

Определение 6.5. Квадратурную формулу

$$S(f) \approx \sum_{i=0}^n A_i f(x_i)$$

будем называть *симметричной*, если она имеет симметричные узлы относительно середины отрезка $[a, b]$,

$$x_i - a = b - x_{n-i}, \quad \forall i = \overline{0, n},$$

а также симметричные коэффициенты:

$$A_i = A_{n-i}, \quad \forall i = \overline{0, n}.$$

Лемма 6.1. Пусть весовая функция ρ является четной относительно середины $[a, b]$, а квадратурная формула $S(f) \approx Q_n(f)$ — симметричной. Тогда для всякой нечетной относительно $(a+b)/2$ функции f ,

$$f(x) = -f(a+b-x),$$

справедливо

$$S(f) = Q_n(f).$$

Доказательство. В силу нечетности f имеем $S(f) = 0$, а симметрия узлов дает $f(x_i) = -f(x_{n-i})$. Покажем, что $Q_n(f) = 0$:

$$Q_n(f) = \sum_{i=0}^n A_i f(x_i) = - \sum_{i=0}^n A_{n-i} f(x_{n-i}) = -Q_n(f),$$

откуда сразу следует утверждение леммы. \square

Теорема 6.3 (о повышенной АСТ симметричных КФ). *Пусть весовая функция ρ является четной относительно середины $[a, b]$, а квадратурная формула $S(f) \approx Q_n(f)$ точна для всех многочленов степени $2M$ и ниже, $M \in \mathbb{Z}$. Если к тому же эта КФ симметрична, то ее алгебраическая степень точности равна как минимум $2M + 1$.*

Доказательство. Нам достаточно показать, что $Q_n(P) = S(P)$ для любого многочлена P степени $2M + 1$. Итак, пусть $\deg P = 2M + 1$. Его старший коэффициент обозначим α и рассмотрим многочлен

$$U(x) = \alpha \left(x - \frac{a+b}{2} \right)^{2M+1}.$$

В силу нечетности этого многочлена $S(U) = 0$, а также в силу симметричности по лемме 6.1 имеем $Q_n(U) = 0$.

Осталось представить P в виде

$$P = U + \tilde{P},$$

где $\deg \tilde{P} = \deg(P - U) \leq 2M$. Тогда по условию $S(\tilde{P}) = Q_n(\tilde{P})$ и

$$S(P) = S(U) + S(\tilde{P}) = S(\tilde{P}) = Q_n(\tilde{P}) = Q_n(\tilde{P}) + Q_n(U) = Q_n(P). \quad \square$$

▷₃ Укажите, в каком месте доказательства используется четность ρ .

6.3.3. Симметрия интерполяционных квадратурных формул

Теорема 6.4. *Пусть функция ρ является четной относительно середины отрезка $[a, b]$. Тогда если узлы $\{x_i\}$ симметричны, то соответствующая им интерполяционная квадратурная формула с весом ρ является симметричной.*

Доказательство. Нам нужно доказать, что если узлы удовлетворяют свойству $x_i - a = b - x_{n-i}$, то коэффициенты соответствующей КФ, которые вычисляются по формуле (6.7):

$$A_i = \int_a^b \Lambda_i(x) \rho(x) dx,$$

обладают свойством $A_i = A_{n-i}$. Доказательство основывается на факте попарной симметричности базисных функций Λ_i :

$$\Lambda_i(x) = \Lambda_{n-i}(a + b - x).$$

Докажем это:

$$\begin{aligned} \Lambda_{n-i}(a + b - x) &= \prod_{j \neq n-i} \frac{a + b - x - x_j}{x_{n-i} - x_j} = \\ &= \prod_{j \neq n-i} \frac{x_{n-j} - x}{(a + b - x_i) - (a + b - x_{n-j})} = \\ &= \prod_{j \neq n-i} \frac{x - x_{n-j}}{x_i - x_{n-j}} = \prod_{j \neq i} \frac{x - x_j}{x_i - x_j} = \Lambda_i(x). \end{aligned}$$

Осталось сделать замену переменных в интеграле:

$$\begin{aligned} A_i &= \int_a^b \Lambda_i(x) \rho(x) dx = \int_a^b \Lambda_{n-i}(a + b - x) \rho(x) dx = [z = a + b - x] = \\ &= - \int_b^a \Lambda_{n-i}(z) \rho(a + b - z) dz = \int_a^b \Lambda_{n-i}(z) \rho(z) dz = A_{n-i}. \quad \square \end{aligned}$$

Следствие 6.1. Интерполяционная квадратурная формула с четной относительно середины $[a, b]$ весовой функцией ρ и симметричными узлами $\{x_i\}$ при $n = 2M$ (т. е. при нечетном количестве узлов) имеет алгебраическую степень точности не менее $2M + 1$.

Доказательство. Поскольку рассматриваемая формула интерполяционная, она точна для всех многочленов степени $n = 2M$ и ниже. По теореме 6.4 она также является симметричной. Значит, по теореме 6.3 АСТ этой формулы равна как минимум $2M + 1$. \square

6.3.4. Квадратурные формулы Ньютона – Котеса

Определение 6.6. Интерполяционные квадратурные формулы для весовой функции $\rho \equiv 1$ по равноотстоящим узлам на отрезке $[a, b]$,

$$x_i = a + \frac{b-a}{n}i, \quad i = \overline{0, n},$$

называются *формулами Ньютона – Котеса*.

По определению и теореме 6.4 все формулы Ньютона – Котеса являются симметричными. С одной такой формулой мы уже знакомы: при $n = 1$ получается формула трапеций. Под формальное определение попадает также формула левых прямоугольников, хотя по смыслу больше подходит формула прямоугольников средних. Сейчас мы построим еще одну классическую квадратурную формулу Ньютона – Котеса.

Формула Симпсона

Квадратурная формула Симпсона — это формула Ньютона – Котеса с тремя узлами:

$$x_0 = a, \quad x_1 = \frac{a+b}{2}, \quad x_2 = b.$$

Иногда эту формулу также называют формулой парабол (подынтегральная функция приближается параболой). Для вычисления коэффициентов этой формулы можно воспользоваться определением (6.7), но мы сделаем это по-другому.

Начнем с того, что выведем формулу Симпсона для частного случая $[a, b] = [-1, 1]$. Мы знаем, что эта формула является симметричной, т. е. имеет вид

$$\int_{-1}^1 f(x)dx = S(f) \approx \bar{Q}_2(f) = \bar{A}_0 f(-1) + \bar{A}_1 f(0) + \bar{A}_0 f(1).$$

Эта формула по построению точна для всех многочленов степени 2 и ниже, а значит при $f(x) = x^i$ для $i = 0, 1, 2$ будет иметь место тождество $S(f) = \bar{Q}_2(f)$. Эти три условия дают систему уравнений

$$\begin{cases} \bar{A}_0 + \bar{A}_1 + \bar{A}_0 = 2, \\ -\bar{A}_0 + \bar{A}_0 = 0, \\ \bar{A}_0 + \bar{A}_0 = \frac{2}{3}, \end{cases}$$

откуда легко находим

$$\bar{A}_0 = \frac{1}{3}, \quad \bar{A}_1 = \frac{4}{3}.$$

Мы получили формулу Симпсона для отрезка $[-1, 1]$:

$$\int_{-1}^1 f(x)dx \approx \frac{1}{3}f(-1) + \frac{4}{3}f(0) + \frac{1}{3}f(1). \quad (6.14)$$

Для вывода формулы в общем случае теперь достаточно сделать замену переменной интегрирования:

$$\int_a^b f(x)dx = \frac{b-a}{2} \int_{-1}^1 f\left(\frac{a+b}{2} + \frac{b-a}{2}t\right)dt.$$

Применяя к последнему интегралу формулу (6.14), окончательно получим *квадратурную формулу Симпсона*

$$\int_a^b f(x)dx \approx \frac{b-a}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right). \quad (6.15)$$

▷₄ Аналогичным образом выведите формулы Ньютона – Котеса с четырьмя и пятью узлами.

▷₅ Определите АСТ формулы Симпсона и используя эту информацию получите выражение для ее остатка путем разложения в ряд Тейлора в точке $(a+b)/2$.

6.4. Квадратурные формулы Гаусса

6.4.1. Постановка задачи

Рассмотрим задачу построения квадратурных формул

$$S(f) = \int_a^b f(x)\rho(x)dx \approx \sum_{i=0}^n A_i f(x_i) = Q_n(f),$$

имеющих максимально возможную алгебраическую степень точности. Такие формулы называются *гауссовыми* или *квадратурными формулами наивысшей АСТ*. Как мы помним, по определению АСТ квадратурной формулы равна m , если

$$S(p) = Q_n(p) \quad \forall p \in \mathbb{P}_m \quad (6.16a)$$

и

$$\exists \tilde{p} \in \mathbb{P}_{m+1} : \quad S(\tilde{p}) \neq Q_n(\tilde{p}). \quad (6.16б)$$

Здесь и далее \mathbb{P}_m — множество всех многочленов степени не выше m . Мы помним также, что в силу линейности для выполнения условия (6.16a) необходимо и достаточно потребовать, чтобы квадратурная формула работала точно на базисе пространства \mathbb{P}_m . Проще всего выбрать степенной базис, поэтому условия (6.16) эквивалентны (в частности) следующим:

$$\sum_{i=0}^n A_i x_i^j = \int_a^b x^j \rho(x) dx, \quad j = \overline{0, m}; \quad (6.17a)$$

$$\sum_{i=0}^n A_i x_i^{m+1} \neq \int_a^b x^{m+1} \rho(x) dx. \quad (6.17б)$$

В дальнейшем мы будем рассматривать только условия (6.17a).

Какова же будет максимально возможная АСТ для данного n ? Условия (6.17) представляют собой систему из $(m+1)$ нелинейных уравнений относительно $(2n+2)$ неизвестных $\{x_i\}_{i=0}^n, \{A_i\}_{i=0}^n$. Поэтому можно по крайней мере надеяться, что эта система будет иметь решение, когда количество уравнений и неизвестных совпадает, т. е. когда

$$m = 2n + 1.$$

Как будет показано далее, это и есть максимально возможная АСТ.

6.4.2. Построение квадратурных формул Гаусса

Итак, построение квадратурных формул Гаусса заключается в нахождении узлов $\{x_i\}$ и коэффициентов $\{A_i\}$, удовлетворяющих условиям (6.17a) для $m = 2n + 1$. Самый очевидный способ их нахождения — это непосредственное решение указанной системы. Для каждого n мы будем иметь $2n + 2$ нелинейных уравнений. Случай $n = 0$ прост:

$$\begin{cases} A_0 = \int_a^b \rho(x) dx, \\ A_0 x_0 = \int_a^b x \rho(x) dx, \end{cases}$$

откуда, в частности, при $\rho(x) \equiv 1$ получим формулу средних прямоугольников. При $n = 1$ систему можно решить относительно легко, а при больших значениях n задача становится практически неподъемной.

▷₁ Решите систему при $n = 1$, $[a, b] = [-1, 1]$ и $\rho(x) \equiv 1$.

Следующее тривиальное следствие теоремы 6.2 существенно упрощает задачу: *если квадратурная формула имеет наивысшую АСТ, то она является интерполяционной, т. е. ее коэффициенты имеют вид*

$$A_i = \int_a^b \prod_{j \neq i} \frac{x - x_j}{x_i - x_j} \rho(x) dx, \quad i = \overline{0, n}.$$

▷₂ Докажите.

Таким образом, задача *упростилась в два раза*: достаточно найти узлы формулы Гаусса, а коэффициенты можно определить по формуле, приведенной выше. О том, какими должны быть узлы формулы, говорит следующая далее теорема. Она является наиболее важным результатом этого пункта, и для ее формулировки понадобится определение скалярного произведения и ортогональности функций.

Определение 6.7. Скалярным произведением интегрируемых с квадратом функций u и v называется выражение

$$(u, v) = \int_a^b u(x)v(x)\rho(x)dx,$$

Здесь ρ — заданная неотрицательная на $[a, b]$ весовая функция. Если $(u, v) = 0$, то говорят, что функции u и v являются *ортогональными* по весу ρ .

Теорема 6.5 (К. Ф. Гаусс). Для того чтобы квадратурная формула с $(n + 1)$ узлами была точной для любых алгебраических многочленов степени $(2n + 1)$ и ниже, необходимо и достаточно, чтобы:

- 1) она была интерполяционной;
- 2) многочлен $\omega_{n+1}(x) = (x - x_0) \dots (x - x_n)$ был ортогонален по весу ρ на отрезке $[a, b]$ ко всем многочленам степени n и ниже:

$$(\omega_{n+1}, p) = \int_a^b \omega_{n+1}(x) p(x) \rho(x) dx = 0 \quad \forall p \in \mathbb{P}_n. \quad (6.18)$$

Доказательство. Пусть формула имеет АСТ $2n + 1$. Значит, она будет точна для всех многочленов вида $q(x) = \omega_{n+1}(x)p(x)$, $p \in \mathbb{P}_n$, так как $\deg q \leq 2n + 1$. С учетом того, что $\omega_{n+1}(x_i) = 0$ при всех i , получим

$$S(q) = (\omega_{n+1}, p) = Q_n(q) = \sum_{i=0}^n A_i \omega_{n+1}(x_i) p(x_i) = 0.$$

То, что формула интерполяционная, было доказано ранее.

Теперь пусть верно (6.18). Возьмем любой $q \in \mathbb{P}_{2n+1}$ и разделим его на ω_{n+1} с остатком:

$$q(x) = \omega_{n+1}(x)p(x) + r(x),$$

где $\deg p \leq n$, $\deg r \leq n$. Тогда с учетом $q(x_i) = r(x_i)$ имеем

$$S(q) = (\omega_{n+1}, p) + S(r) = S(r) = \sum_{i=0}^n A_i r(x_i) = \sum_{i=0}^n A_i q(x_i),$$

т. е. формула точна для любого $q \in \mathbb{P}_{2n+1}$. □

Таким образом, узлами квадратурной формулы Гаусса являются корни многочлена, ортогонального по весу ρ всем многочленам степени n и ниже. Многочлен, удовлетворяющий этим условиям, в дальнейшем будем обозначать ω_{n+1}^* .

Для строгости доказательства полученного результата нужно еще: а) показать, что нельзя достичь АСТ, большей $2n + 1$; б) доказать существование узлов формул Гаусса на отрезке $[a, b]$ (вдруг они вообще комплексные?).

Теорема 6.6. Если весовая функция ρ знакопостоянна на отрезке $[a, b]$, то не существует квадратурных формул с $(n + 1)$ узлами, имеющих АСТ $2n + 2$.

Доказательство. Предположим существование такой формулы. Пусть $\{x_i\}_{i=0}^n$ — ее узлы. Рассмотрим многочлен

$$\mathbf{v}(x) = (x - x_0)^2 \dots (x - x_n)^2.$$

Для него очевидно имеем $S(\mathbf{v}) \neq 0$, так как подынтегральное выражение знакопостоянно и не равно нулю тождественно. С другой стороны, имеем $Q_n(f) = \sum_{i=0}^n A_i \mathbf{v}(x_i) = 0$. Полученное противоречие доказывает теорему. \square

Теорема 6.7. *Если весовая функция ρ сохраняет знак на отрезке $[a, b]$, то многочлен ω_{n+1}^* степени $(n + 1)$, ортогональный на данном отрезке по весу ρ ко всем многочленам степени n и ниже, существует и единствен для любого фиксированного n . При этом все его корни действительны, различны и лежат внутри $[a, b]$.*

Замечание 6.2. Рассмотрим множество многочленов $\{\omega_i^*\}_{i=1}^\infty$. Поскольку по определению $(\omega_{n+1}^*, p) = 0$ для всех $p \in \mathbb{P}_n$, то имеем

$$(\omega_i^*, \omega_j^*) = 0, \quad \text{если } i \neq j,$$

т. е. множество $\{\omega_i^*\}$ представляет собой *систему ортогональных многочленов* (по весу ρ). Такие системы играют важную роль во многих приложениях. Поэтому результат главной теоремы 6.5 можно сформулировать и так: *квадратурная формула с $(n + 1)$ узлами точна для любых алгебраических многочленов степени $(2n + 1)$ и ниже тогда и только тогда, когда она интерполяционная, а ее узлы являются корнями ортогонального по весу ρ многочлена степени $n + 1$.*

6.4.3. Классические квадратурные формулы Гаусса

Единичная весовая функция

Рассмотрим случай единичной весовой функции на отрезке $[-1, 1]$, т. е. квадратурные формулы наивысшей АСТ для интеграла общего вида $\int_a^b f(x)dx$. Такие формулы иногда называют формулами Гаусса — Лежандра, так как их узлами являются корни ортогональных на $[-1, 1]$ с единичным весом *многочленов Лежандра*. Эти многочлены можно определить явной формулой

$$L_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} ((x^2 - 1)^n) \quad (6.19)$$

или рекуррентным соотношением

$$\begin{aligned} L_{n+1}(x) &= \frac{2n+1}{n+1}xL_n(x) - \frac{n}{n+1}L_{n-1}(x), \\ L_0(x) &= 1, \quad L_1(x) = x. \end{aligned} \quad (6.20)$$

▷₃ Постройте квадратурные формулы Гаусса – Лежандра с двумя и тремя узлами.

Для применения квадратурных формул Гаусса – Лежандра на произвольном отрезке $[a, b]$ достаточно сделать уже известную нам замену переменных

$$x = \frac{a+b}{2} + \frac{b-a}{2}t, \quad t \in [-1, 1],$$

как это было при выводе формулы Симпсона. Покажем подробно, как это делается. Пусть имеется формула Гаусса вида

$$\int_{-1}^1 f(x)dx \approx \sum_{i=0}^n A_i f(x_i).$$

Тогда

$$\begin{aligned} \int_a^b f(x)dx &= \frac{b-a}{2} \int_{-1}^1 \underbrace{f\left(\frac{a+b}{2} + \frac{b-a}{2}t\right)}_{g(t)} dt = \frac{b-a}{2} \int_{-1}^1 g(t)dt \approx \\ &\approx \frac{b-a}{2} \sum_{i=0}^n A_i g(x_i) = \frac{b-a}{2} \sum_{i=0}^n A_i f\left(\frac{a+b}{2} + \frac{b-a}{2}x_i\right). \end{aligned}$$

Другие классические формулы

Для некоторых весовых функций известны явные формулы, определяющие ортогональные многочлены, а также явные рекуррентные соотношения. Для этих многочленов давно составлены таблицы коэффициентов и узлов соответствующих формул наивысшей АСТ. Эти коэффициенты выводить слишком долго, а запоминать бессмысленно, но нужно о них знать, чтобы при необходимости найти в литературе [7] (см. таблицу).

В частности, знакомые нам многочлены Чебышева

$$T_n(x) = \cos(n \arccos x),$$

Классические системы ортогональных многочленов

$\rho(x)$	Отрезок	Название
$\frac{1}{\sqrt{1-x^2}}$	$[-1, 1]$	Многочлены Чебышева
$x^\alpha e^{-x}$	$[0, +\infty)$	Многочлены Лагерра
e^{-x^2}	$(-\infty, +\infty)$	Многочлены Эрмита
$(1-x)^\alpha(1+x)^\beta$	$[-1, 1]$	Многочлены Якоби

оказывается, ортогональны на отрезке $[-1, 1]$ с весом $\rho(x) = (1-x^2)^{-1/2}$. Следовательно, узлы квадратурной формулы Гаусса — Чебышева, которая имеет наивысшую АСТ для интегралов с этой весовой функцией, являются корнями многочленов T_n (чебышевскими узлами):

$$x_i = \cos \frac{\pi(2i+1)}{2n+2}, \quad i = \overline{0, n}.$$

Замечателен тот факт, что коэффициенты таких квадратурных формул при данном n будут равны между собой:

$$A_0 = \dots = A_n = \frac{\pi}{n+1}.$$

Другие весовые функции

Возникает вопрос: каким образом строить гауссовы квадратурные формулы в общем случае, когда неизвестны соответствующие системы ортогональных многочленов? Ответ таков: нужные системы ортогональных многочленов можно построить по специальным рекуррентным соотношениям, рассмотрение которых выходит за рамки нашего курса.

Однако при малых n можно использовать следующий простой подход. Он позволяет для данного n найти непосредственно многочлен ω_{n+1}^* в явном виде

$$\omega_{n+1}^*(x) = x^{n+1} + c_n x^n + c_{n-1} x^{n-1} + \dots + c_1 x + c_0. \quad (6.21)$$

Воспользуемся определением: ω_{n+1}^* должен быть ортогонален всем многочленам степени n и ниже, а для этого необходимо и достаточно потре-

бывать его ортогональности одночленам $\chi_i(x) = x^i$:

$$\int_a^b \omega_{n+1}^*(x) x^i \rho(x) dx = 0, \quad i = \overline{0, n}.$$

Подставляя сюда представление (6.21), получим СЛАУ вида

$$\gamma_{i0}c_0 + \gamma_{i1}c_1 + \dots + \gamma_{i,n-1}c_{n-1} + \gamma_{i,n}c_n = -g_i, \quad i = \overline{0, n}, \quad (6.22)$$

где

$$\gamma_{ij} = \int_a^b x^{i+j} \rho(x) dx, \quad g_i = \int_a^b x^{n+i+1} \rho(x) dx.$$

▷₄ Выведите эти формулы.

После того как найден многочлен ω_{n+1}^* , остается найти его корни и вычислить коэффициенты квадратурной формулы по правилу (6.7).

6.5. Составные квадратурные формулы

6.5.1. Введение

Пусть $\rho \equiv 1$. Рассмотрим задачу вычисления интеграла

$$S(f) = \int_a^b f(x) dx$$

с некоторой заданной точностью. Для этого разобьем отрезок $[a, b]$ на N равных частей точками $\{\xi_k\}_{k=0}^N$,

$$\xi_k = a + kh, \quad k = \overline{0, N}, \quad h = \frac{b - a}{N},$$

и получим

$$S(f) = \sum_{k=1}^N S_k(f), \quad S_k(f) = \int_{\xi_{k-1}}^{\xi_k} f(x) dx.$$

Приближая каждое слагаемое какой-нибудь простой квадратурной формулой

$$S_k(f) \approx Q_{n,k}(f),$$

получим *составную квадратурную формулу* для всего интеграла, которую будем обозначать как

$$S(f) \approx Q_n^N(f) = \sum_{k=1}^N Q_{n,k}(f).$$

6.5.2. Простейшие составные формулы

Возьмем в качестве базовой, например, формулу средних прямоугольников (6.2):

$$\int_{\xi_{k-1}}^{\xi_k} f(x) dx \approx (\xi_k - \xi_{k-1}) f\left(\frac{\xi_{k-1} + \xi_k}{2}\right).$$

С учетом того, что $\xi_k = a + kh$, после суммирования всех частей получим *составную формулу средних прямоугольников* (рис. 6.1):

$$\int_a^b f(x)dx \approx h \sum_{k=1}^N f\left(a + \left(k - \frac{1}{2}\right)h\right). \quad (6.23)$$

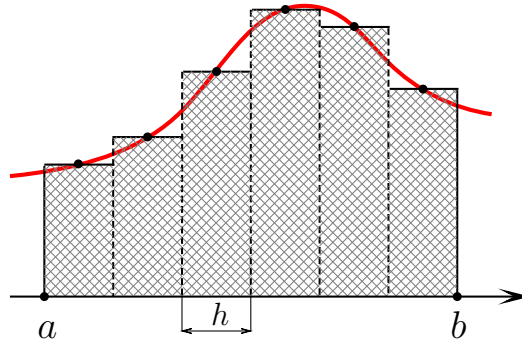


Рис. 6.1. Составная квадратурная формула средних прямоугольников

Аналогично строится *составная формула трапеций* (рис. 6.2):

$$\int_{\xi_{k-1}}^{\xi_k} f(x)dx \approx \frac{h}{2}(f(\xi_{k-1}) + f(\xi_k)),$$

откуда

$$\int_a^b f(x)dx \approx \frac{h}{2}\left(f(a) + 2 \sum_{k=1}^{N-1} f(a + kh) + f(b)\right). \quad (6.24)$$

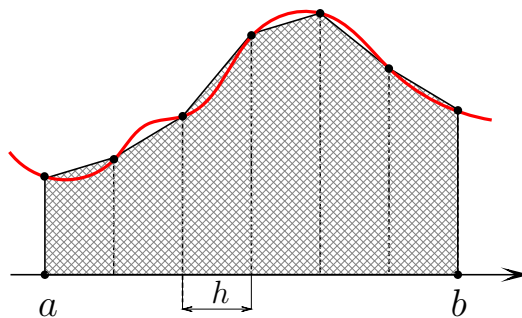


Рис. 6.2. Составная квадратурная формула трапеций

▷₁ Постройте составные формулы левых, правых прямоугольников и Симпсона.

6.5.3. Остаток составных квадратурных формул

Теперь наша задача — вычислить остаток составных формул, который обозначим

$$R_n^N(f) = S(f) - Q_n^N(f).$$

Согласно введенным в начале пункта обозначениям, имеем

$$\begin{aligned} R_n^N(f) &= S(f) - Q_n^N(f) = \sum_{k=1}^N S_k(f) - \sum_{k=1}^N Q_{n,k}(f) = \\ &= \sum_{k=1}^N (S_k(f) - Q_{n,k}(f)) = \sum_{k=1}^N R_{n,k}(f). \end{aligned}$$

Таким образом, как нетрудно было догадаться сразу, *остаток составной квадратурной формулы равен сумме остатков ее составных частей*.

Предположим, что для каждой квадратурной суммы $Q_{n,k}$ имеет место представление остатка в виде

$$R_{n,k}(f) = Ch^l f^{(m)}(\eta_k),$$

где $\eta_k \in [\xi_{k-1}, \xi_k]$; C — известная константа; l и m — фиксированные натуральные числа. Как мы помним, такая форма остатка справедлива, в частности, для всех уже рассмотренных нами интерполяционных формул (6.11), (6.12), (6.13). Тогда для всего остатка получаем очевидное представление

$$R_n^N(f) = \sum_{k=1}^N R_{n,k}(f) = Ch^l \sum_{k=1}^N f^{(m)}(\eta_k).$$

Эту формулу можно упростить, если предположить, что $f^{(m)}$ непрерывна на $[a, b]$. В этом случае на этом отрезке она достигает своих экстремальных значений y_{\min} и y_{\max} . Тогда

$$y_{\min} \leq \frac{1}{N} \sum_{k=1}^N f^{(m)}(\eta_k) \leq y_{\max}.$$

По теореме о промежуточном значении, в силу непрерывности $f^{(m)}$ на отрезке $[a, b]$ эта функция «проходит» все значения от y_{\min} до y_{\max} , а в

частности, существует такая точка $\eta \in [a, b]$, что

$$f^{(m)}(\eta) = \frac{1}{N} \sum_{k=1}^N f^{(m)}(\eta_k). \quad (6.25)$$

Следовательно, формула для остатка составной КФ примет вид

$$R_n^N(f) = Ch^l N f^{(m)}(\eta) \quad (6.26a)$$

или, с учетом $N = (b - a)/h$,

$$R_n^N(f) = C(b - a)h^{l-1} f^{(m)}(\eta). \quad (6.26a)$$

Замечание 6.3. Напомним, что здесь η — некоторая точка из отрезка $[a, b]$, причем эта точка зависит от разбиения, т. е. для различных N эти точки будут разными (η зависит от h).

Отсюда для составной формулы средних прямоугольников (6.23) из (6.12) получаем остаток

$$R_0^N(f) = \frac{b-a}{24} h^2 f''(\eta), \quad (6.27)$$

а для составной формулы трапеций (6.24) из (6.13) имеем

$$R_1^N(f) = -\frac{b-a}{12} h^2 f''(\eta). \quad (6.28)$$

Имея формулы такого рода, теоретически можно заранее оценить необходимое число отрезков разбиения N , которое следует взять для вычисления интеграла с требуемой точностью. Для этого, естественно, нужно располагать информацией о функции $f^{(m)}$ на отрезке $[a, b]$. На практике такая информация доступна не всегда, поэтому в универсальных программах численного интегрирования используется несколько иной способ представления погрешности составных формул.

6.5.4. Выделение главной части погрешности

Рассмотрим формулу остатка составной формулы прямоугольников (6.27). Она говорит, что остаток имеет вид $O(h^2)$, если f'' существует и непрерывна на $[a, b]$:

$$S(f) = Q_0^N(f) + O(h^2). \quad (6.29)$$

Если же f имеет производные более высоких порядков, то из этого остатка можно выделить *главную часть*, которая, в отличие от $f''(\eta)$, не зависит от h и может быть вычислена. Покажем, как это делается на примере той же составной формулы средних прямоугольников (6.23).

Рассмотрим точку $z_k = (\xi_{k-1} + \xi_k)/2$ — середину отрезка $[\xi_{k-1}, \xi_k]$. Предположим, что f четырежды непрерывно дифференцируема и представим ее формулой Тейлора в точке z_k :

$$S_k(f) = \int_{\xi_{k-1}}^{\xi_k} f(x) dx = \int_{\xi_{k-1}}^{\xi_k} \left[f(z_k) + (x - z_k)f'(z_k) + \frac{(x - z_k)^2}{2}f''(z_k) + \frac{(x - z_k)^3}{6}f'''(z_k) + \frac{(x - z_k)^4}{24}f^{(4)}(\zeta_k) \right] dx.$$

Здесь ζ_k — некоторая точка из $[\xi_{k-1}, \xi_k]$. Так как интеграл от нечетных степеней $(x - z_k)$ равен нулю, после интегрирования каждого слагаемого и применения теоремы о среднем к последнему из них получим

$$S_k(f) = hf(z_k) + \frac{h^3}{24}f''(z_k) + \frac{h^5}{1920}f^{(4)}(\eta_k),$$

η_k — точка из $[\xi_{k-1}, \xi_k]$, отличная, вообще говоря, от ζ_k . Просуммируем теперь все $S_k(f)$:

$$S(f) = \sum_{k=1}^N S_k(f) = h \sum_{k=1}^N f(z_k) + \frac{h^3}{24} \sum_{k=1}^N f''(z_k) + \frac{h^5}{1920} \sum_{k=1}^N f^{(4)}(\eta_k).$$

Итак, наш интеграл состоит из трех слагаемых. Рассмотрим отдельно каждое из них. Первое, $h \sum_{k=1}^N f(z_k) = Q_0^N(f)$, — это как раз составная квадратурная сумма средних прямоугольников (6.23). Второе,

$$\frac{h^3}{24} \sum_{k=1}^N f''(z_k) = \frac{h^2}{24} \left(h \sum_{k=1}^N f''(z_k) \right),$$

содержит составную сумму средних прямоугольников для f'' , так что с учетом (6.29) это слагаемое примет вид

$$\frac{h^3}{24} \sum_{k=1}^N f''(z_k) = \frac{h^2}{24} \left(\int_a^b f''(x) dx + O(h^2) \right).$$

В третьем слагаемом применим формулу (6.25) и получим

$$\frac{h^5}{1920} \sum_{k=1}^N f^{(4)}(\eta_k) = \frac{h^5}{1920} N f^{(4)}(\eta) = \frac{h^4}{1920} (b-a) f^{(4)}(\eta).$$

Собирая все вместе, имеем

$$S(f) = Q_0^N(f) + \frac{h^2}{24} \int_a^b f''(x) dx + O(h^4),$$

или

$$R_0^N(f) = \frac{f'(a) - f'(b)}{24} h^2 + O(h^4). \quad (6.30)$$

Таким образом, мы выделили главную часть остатка составной квадратной формулы средних прямоугольников: она равна

$$\frac{f'(a) - f'(b)}{24} h^2.$$

Еще раз обратим внимание, что константа при h^2 , как и было обещано, зависит только от f и, в принципе, может быть вычислена. Эту главную часть при желании можно добавить к исходной квадратурной сумме и получить уточненное значение интеграла $S(f)$. Аналогичную операцию выделения главной части можно провести и с другими составными квадратурными формулами.

6.5.5. Практическая оценка погрешности: правило Рунге

Составные квадратурные формулы по построению представляют собой инструмент, позволяющий регулировать точность вычисления интеграла за счет увеличения числа отрезков разбиения N (уменьшения шага h). Для вычислений на практике нужно уметь выбирать такое N , которое бы позволило достичь требуемой точности ε , и в то же время не приводило к избыточным вычислениям (не было слишком велико). Один из возможных способов реализации этих требований называется *правилом Рунге* или *методом двойного пересчета*.

Пусть $Q_n^N = Q^N$ — некоторая составная квадратурная сумма. Предположим, что для остатка соответствующей квадратурной формулы справедливо представление, аналогичное (6.30):

$$R^N(f) = S(f) - Q^N(f) = Kh^p + o(h^p), \quad h = \frac{b-a}{N}, \quad (6.31)$$

где константа K не зависит от h . Здесь Kh^p представляет собой главную часть погрешности, которую можно найти, если известно значение константы K .

Вычислим приближенное значение интеграла $S(f)$ по данной формуле *дважды*, на двух разных разбиениях с числом отрезков N_1 и N_2 , $N_2 > N_1$. Полученные приближения обозначим соответственно

$$q_1 = Q^{N_1}(f) \quad \text{и} \quad q_2 = Q^{N_2}(f).$$

Обозначив $h_i = (b-a)/N_i$, согласно (6.31) можно записать приближенные равенства

$$S(f) - q_1 \approx Kh_1^p,$$

$$S(f) - q_2 \approx Kh_2^p.$$

Отсюда, исключая $S(f)$, имеем

$$K \approx \frac{q_2 - q_1}{h_1^p - h_2^p} =: \tilde{K},$$

а это очень важная информация. Во-первых, с ее помощью можно получить оценку погрешности для обеих квадратурных сумм:

$$R^{N_i}(f) \approx Kh_i^p \approx \tilde{K}h_i^p = \tilde{R}_i,$$

где

$$\tilde{R}_i = \frac{(q_2 - q_1)h_i^p}{h_1^p - h_2^p}. \quad (6.32)$$

Во-вторых, имея \tilde{R}_i , мы можем уточнить значения q_i :

$$S(f) = q_i + R^{N_i}(f) \approx q_i + \tilde{R}_i. \quad (6.33)$$

Для получения максимальной точности лучше всего уточнять q_2 .

В-третьих, с помощью \tilde{K} можно оценить оптимальную величину шага h^* для достижения точности ε . Делается это с помощью очевидного требования

$$|R^N(f)| \approx |\tilde{K}|h^p \leq \varepsilon,$$

из которого получим

$$h \leq \left(\frac{\varepsilon}{|\tilde{K}|} \right)^{1/p} = \left(\frac{\varepsilon}{|\tilde{R}_i|} \right)^{1/p} h_i = h^*,$$

или

$$N \geq N^* = \frac{b-a}{h^*} = \left(\frac{|\tilde{R}_i|}{\varepsilon} \right)^{1/p} N_i. \quad (6.34)$$

Практические рекомендации. Теперь поговорим о том, как все это реализуется на практике. Простейший алгоритм выглядит следующим образом:

- 1) выбираем квадратурную формулу Q^N , N_1 и N_2 , точность ε ;
- 2) $q_1 \leftarrow Q^{N_1}(f)$, $q_2 \leftarrow Q^{N_2}(f)$;
- 3) по формуле (6.32) вычисляем оценку погрешности для приближения q_2 :

$$\tilde{R} \leftarrow \frac{(q_2 - q_1)h_2^p}{h_1^p - h_2^p}.$$

Как мы помним, здесь p — показатель точности КФ, определяемый (6.31), $h_i = (b-a)/N_i$;

- 4) если $|\tilde{R}| \leq \varepsilon$, полагаем $q \leftarrow q_2$ и завершаем алгоритм. В противном случае переходим к следующему шагу;

- 5) вычисляем «оптимальное» количество разбиений N^* по формуле (6.34):

$$N^* \leftarrow \left\lceil \left(\frac{|\tilde{R}|}{\varepsilon} \right)^{1/p} N_2 \right\rceil.$$

Здесь $\lceil x \rceil$ — минимальное целое, большее x ;

- 6) полагаем $N_1 \leftarrow N_2$, $q_1 \leftarrow q_2$, $N_2 \leftarrow N^*$, $q_2 \leftarrow Q^{N^*}(f)$;
- 7) переходим к шагу 3.

Результатом работы алгоритма является значение $q \approx S(f)$. Этот алгоритм можно совершенствовать в следующих направлениях:

- **экономия вычислений f :** в зависимости от используемой базовой КФ подбирать N_1 и N_2 так, чтобы в соответствующих сетках для вычисления q_1 и q_2 было как можно больше одинаковых узлов. Например, в случае составной формулы трапеций (6.24) обычно берут $N_{i+1} = 2N_i$. Тогда получается

$$Q_1^N(f) = \frac{h}{2} \left(f(a) + 2 \sum_{k=1}^{N-1} f(a + kh) + f(b) \right),$$

$$Q_1^{2N}(f) = \frac{h}{4} \left(f(a) + 2 \sum_{k=1}^{2N-1} f(a + k\frac{h}{2}) + f(b) \right) =$$

$$= \frac{1}{2} \left(Q_1^N(f) + h \sum_{k=1}^N f(a + (2k-1)\frac{h}{2}) \right);$$

▷₂ Постройте соответствующую модификацию алгоритма.

- **использование уточненного приближения:** на шаге 4 можно возвращать не q_2 , а его уточненное по формуле (6.33) значение $q_2 + \tilde{R}$;

- **использование неравномерных сеток** — наиболее радикальное изменение, которое по сути представляет собой совершенно иной алгоритм. Мотивация тут следующая: численное интегрирование по равномерно расположенным узлам является неоптимальным, если функция f на отрезке интегрирования имеет участки с различным поведением (например, сначала сильно осциллирует, а потом изменяется очень медленно). Поэтому современные программы численного интегрирования используют алгоритмы с неравномерно расположенными узлами. При этом принцип оценки погрешности остается прежним.

Глава 7

Численное решение обыкновенных дифференциальных уравнений

- 7.1. Одношаговые методы решения задачи Коши
- 7.2. Методы Рунге – Кутты
- 7.3. Выбор шага численного интегрирования ОДУ
- 7.4. Многошаговые методы
- 7.5. Численное решение краевых задач. Метод стрельбы
- 7.6. Решение краевых задач конечно-разностным методом

$$y(x) = y_0 + \int_{x_0}^x f(z, y(z)) dz. \quad (7.2)$$

Пусть необходимо найти значение решения в некоторой точке $x_1 > x_0$. На практике редко удастся вычислить $y(x_1)$ точно, поэтому приходится прибегать к приближенным (численным) методам.

7.1.2. Методы Эйлера

Явный метод Эйлера

Рассмотрим скалярный случай $n = 1$. Начальное условие (7.16) позволяет составить уравнение касательной к точному решению y в точке x_0 : если (x_1, y_1) — точка на касательной, то справедливо соотношение

$$\frac{y_1 - y_0}{x_1 - x_0} = y'(x_0) = f(x_0, y_0).$$

Полагая $y(x_1) \approx y_1$, получим самый простой и самый известный метод — *явный метод Эйлера*

$$y_1 = y_0 + (x_1 - x_0)f(x_0, y_0). \quad (7.3)$$

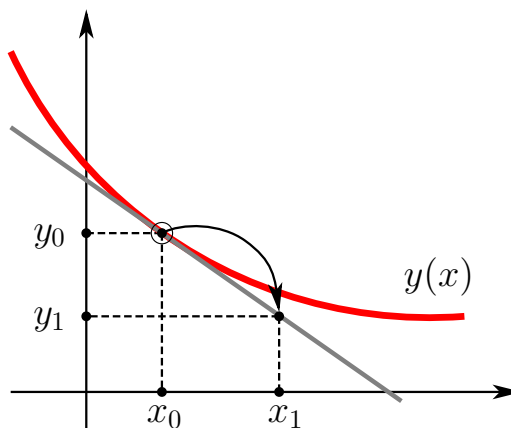


Рис. 7.1.

Геометрический смысл этого метода представлен на рис. 7.1.

Неявный метод Эйлера

Рассмотрим соотношение (7.2) для $x = x_1$:

$$y(x_1) = y_0 + \int_{x_0}^{x_1} f(z, y(z)) dz.$$

Приблизим интеграл по правилу правых прямоугольников⁷

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

и получим неявное соотношение для вычисления $y_1 \approx y(x_1)$:

$$y_1 = y_0 + (x_1 - x_0) f(x_1, y_1). \quad (7.4)$$

Это *неявный* метод Эйлера. Очевидно, что для вычисления y_1 необходимо решать нелинейное уравнение (систему из n нелинейных уравнений в общем случае).

▷₁ В чем состоит геометрический смысл неявного метода Эйлера? Изобразите соответствующий рисунок.

▷₂ Примените для аппроксимации интеграла (7.2) квадратурную формулу трапеций и запишите формулу для соответствующего метода.

7.1.3. Одношаговые методы: терминология

Итак, методы Эйлера позволяют вычислить приближенное значение решения в точке x_1 по известному значению y_0 в точке x_0 . Методы такого типа называются *одношаговыми*.

Определение 7.1. *Одношаговым методом* численного интегрирования задачи Коши (7.1) называется отображение

$$\Phi : \{x_0, y_0, x_1\} \mapsto y_1,$$

которое данному начальному условию (x_0, y_0) ставит в соответствие значение приближенного решения задачи (7.1) в данной точке $x_1 \in \mathbb{R}$:

$$\Phi(x_0, y_0, x_1) = y_1 \approx y(x_1).$$

⁷Заметим, что если применить формулу левых прямоугольников, получится явный метод Эйлера (7.3).

Главное различие между двумя методами Эйлера состоит в том, что в одном случае y_1 задается явной формулой, а в другом оно определяется как решение нелинейного уравнения. Аналогично все методы численного решения ОДУ разделяются на явные и неявные.

Определение 7.2. Если отображение Φ является неявной функцией, т. е. задано уравнением вида

$$\Psi(x_0, y_0, x_1, y_1) = 0,$$

то соответствующий одношаговый метод называется *неявным*. В противном случае, т. е. если y_1 явно выражается через x_0, y_0 и x_1 , метод называется *явным*.

Как правило, отображение Φ определено для всех x_1 , лежащих в некоторой окрестности x_0 , что позволяет рассматривать $\Phi(x_0, y_0, x)$ как функцию непрерывного аргумента x . Поэтому иногда для краткости будем рассматривать отображение $\varphi : \mathbb{R} \rightarrow \mathbb{R}^n$, определенное по правилу

$$\varphi(x) = \Phi(x_0, y_0, x).$$

Обычно вычислителю необходимо получить приближенное решение задачи Коши (7.1) на «большом» отрезке $[x_0, x_0 + H]$. Поскольку функция φ дает приемлемое приближение лишь в достаточно малой окрестности точки x_0 , отрезок интегрирования разбивается на N частей сеткой узлов

$$x_0 < x_1 < x_2 \dots < x_N = x_0 + H \quad (7.5)$$

и строится набор приближенных значений $y_k \approx y(x_k)$ по правилу

$$y_{k+1} = \Phi(x_k, y_k, x_{k+1}), \quad k = \overline{0, N-1}. \quad (7.6)$$

7.1.4. Порядок метода

Классическим показателем точности методов численного интегрирования ОДУ является порядок метода. Это понятие связано с разложением точного и приближенного решений в ряд Тейлора.

Рассмотрим сначала y — точное решение уравнения (7.1) в скалярном случае и предположим, что оно достаточно гладкое. Введем следующие

традиционные обозначения:

$$f = f(x_0, y_0), \quad f_x = \frac{\partial f}{\partial x}(x_0, y_0), \quad f_y = \frac{\partial f}{\partial y}(x_0, y_0), \quad (7.7)$$

$$f \underbrace{x \dots x}_M \underbrace{y \dots y}_N = \frac{\partial^{M+N} f}{\partial x^M \partial y^N}(x_0, y_0). \quad (7.8)$$

Разложение Тейлора для $y(x_0 + h)$ в точке x_0 имеет вид

$$y(x) = y_0 + hy'(x_0) + \frac{h^2}{2!}y''(x_0) + \frac{h^3}{3!}y'''(x_0) + \dots \quad (7.9)$$

Тот факт, что y удовлетворяет (7.1), позволяет вычислить неизвестные коэффициенты $y^{(k)}(x_0)$ для любого k :

$$\begin{aligned} y'(x_0) &= f, \\ y''(x_0) &= \left. \frac{d}{dx} f(x, y(x)) \right|_{x=x_0} = f_x + f_y y'(x_0) = f_x + f_y f, \\ y'''(x_0) &= \left. \frac{d^2}{dx^2} f(x, y(x)) \right|_{x=x_0} = f_{xx} + 2f_{xy}f + f_{yy}f^2 + f_y(f_x + f_y f), \end{aligned} \quad (7.10)$$

и т. д. Таким образом, мы имеем теоретическую возможность *точно* вычислить разложение Тейлора (7.9).

Теперь зафиксируем начальное условие (x_0, y_0) , рассмотрим произвольный одношаговый метод Φ и $\varphi(x) = \Phi(x_0, y_0, x)$. Традиционный способ оценки точности такого приближения заключается в сравнении разложения Тейлора для $\varphi(x)$ с разложением (7.9). Чем больше членов в этих разложениях совпадают, тем выше порядок метода.

Определение 7.3. *Локальной погрешностью* одношагового метода Φ называется функция

$$r(x) = y(x) - \varphi(x) = y(x) - \Phi(x_0, y_0, x).$$

Определение 7.4. Одношаговый метод имеет порядок p , если для всех достаточно гладких задач его локальная погрешность представима в виде

$$r(x_0 + h) = Ch^{p+1} + O(h^{p+2}), \quad C \neq 0, \quad (7.11)$$

т. е. разложения Тейлора в точке x_0 для точного решения $y(x_0 + h)$ и приближенного $\varphi(x_0 + h)$ совпадают до члена h^p включительно:

$$y^{(k)}(x_0) = \varphi^{(k)}(x_0) \quad \forall k = \overline{0, p}.$$

Слагаемое Ch^{p+1} называют *главным членом локальной погрешности*, а константу C — *константой погрешности* метода.

Замечание 7.1. Следует понимать, что высокий порядок метода не всегда гарантирует высокую точность приближения.

Порядок явного метода Эйлера. Согласно определению (7.3), имеем

$$\varphi(x) = \Phi(x_0, y_0, x) = y_0 + (x - x_0)f(x_0, y_0).$$

Отсюда получим

$$\begin{aligned}\varphi(x_0) &= y_0, \\ \varphi'(x_0) &= f(x_0, y_0) = y'(x_0), \\ \varphi^{(k)}(x_0) &= 0 \neq y^{(k)}(x_0) \quad \forall k \geq 2.\end{aligned}$$

Следовательно,

$$r(x_0 + h) = y(x_0 + h) - \varphi(x_0 + h) = \frac{h^2}{2!}y''(x_0) + O(h^3).$$

Таким образом, порядок явного метода Эйлера равен единице, а константа погрешности согласно (7.10) равна

$$C = \frac{1}{2}(f_x + f_y f). \quad (7.12)$$

Порядок неявного метода Эйлера. Рассмотрим (7.4). Для этого метода приближенное решение $\varphi(x)$ определяется неявным соотношением

$$\varphi(x) = y_0 + (x - x_0)f(x, \varphi(x)).$$

Вычислим производные $\varphi(x)$:

$$\begin{aligned}\varphi'(x) &= f(x, \varphi(x)) + (x - x_0)\frac{d}{dx}f(x, \varphi(x)), \\ \varphi''(x) &= 2\frac{d}{dx}f(x, \varphi(x)) + (x - x_0)\frac{d^2}{dx^2}f(x, \varphi(x)).\end{aligned}$$

Отсюда с учетом того, что $\varphi(x_0) = y_0$, получим

$$\begin{aligned}\varphi'(x_0) &= f(x_0, y_0) = y'(x_0), \\ \varphi''(x_0) &= 2\frac{d}{dx}f(x_0, y_0) = 2(f_x + f_y f) \neq y''(x_0).\end{aligned}$$

Следовательно, разложение локальной погрешности в ряд Тейлора имеет вид

$$r(x_0 + h) = y(x_0 + h) - \varphi(x_0 + h) = -\frac{h^2}{2!}y''(x_0) + O(h^3).$$

Итак, неявный метод Эйлера имеет первый порядок, а его константа погрешности равна

$$C = -\frac{1}{2}(f_x + f_y f). \quad (7.13)$$

Замечание 7.2. На рассмотренных примерах видно, что в общем случае константа погрешности C представляет собой вектор, который выражается через значения частных производных f в точке (x_0, y_0) .

▷₃ Найдите порядок метода, построенного в ▷₂, а также главный член его локальной погрешности.

7.2. Методы Рунге — Кутты

Методы типа Рунге — Кутты (РК) являются наиболее популярными одношаговыми методами численного решения обыкновенных дифференциальных уравнений в настоящее время. Название этих методов связано с именами немецких математиков Карла Рунге (1856—1927) и Мартина Кутты (1867—1944).

7.2.1. Простейшие методы Рунге — Кутты

До сих пор мы рассматривали лишь два простейших метода численного решения ОДУ — явный и неявный методы Эйлера. Рассмотрим теперь более совершенные одношаговые методы.

Рассмотрим задачу Коши в форме интегрального уравнения (7.2) для $x = x_0 + h$:

$$y(x_0 + h) = y_0 + \int_{x_0}^{x_0+h} f(z, y(z)) dz.$$

В интеграле удобно сделать замену переменных $z = x_0 + th$:

$$y(x_0 + h) = y_0 + h \int_0^1 f(x_0 + th, y(x_0 + th)) dt. \quad (7.14)$$

Приблизим интеграл квадратурной формулой средних прямоугольников:

$$y(x_0 + h) \approx y_0 + hf\left(x_0 + \frac{h}{2}, y\left(x_0 + \frac{h}{2}\right)\right).$$

Использовать эту формулу для вычислений пока нельзя, так как неизвестно значение $y(x_0 + h/2)$. Но его можно приблизить любым из методов Эйлера (7.3), (7.4). Для начала возьмем явный метод:

$$y\left(x_0 + \frac{h}{2}\right) \approx y_0 + \frac{h}{2}f(x_0, y_0),$$

и получим численный метод вида

$$Y_2 = y_0 + \frac{h}{2}f(x_0, y_0), \quad (7.15a)$$

$$y_1 = y_0 + hf\left(x_0 + \frac{h}{2}, Y_2\right). \quad (7.15б)$$

Этот метод будем называть *явным методом средних прямоугольников (средней точки)*.

Если же приблизить $y\left(x_0 + \frac{h}{2}\right)$ неявным методом Эйлера

$$y\left(x_0 + \frac{h}{2}\right) \approx Y_1 = y_0 + \frac{h}{2}f\left(x_0 + \frac{h}{2}, Y_1\right),$$

получим *неявный метод средних прямоугольников*:

$$Y_1 = y_0 + \frac{h}{2}f\left(x_0 + \frac{h}{2}, Y_1\right), \quad (7.16a)$$

$$y_1 = y_0 + hf\left(x_0 + \frac{h}{2}, Y_1\right). \quad (7.16б)$$

Оба полученных метода, а также и сами методы Эйлера являются частными случаями методов Рунге — Кутты.

7.2.2. Общий случай

В общем случае приблизим интеграл в (7.14) какой-то абстрактной квадратурной формулой с узлами $\{c_1, \dots, c_s\}$ и весами $\{b_1, \dots, b_s\}$:

$$\int_0^1 f(x_0 + th, y(x_0 + th))dt \approx \sum_{i=1}^s b_i f(x_0 + c_i h, y(x_0 + c_i h)).$$

Для нахождения приближений к неизвестным величинам $y(x_0 + c_i h)$ используем такой же подход:

$$y(x_0 + c_i h) = y_0 + \int_{x_0}^{x_0 + c_i h} f(z, y(z))dz = y_0 + h \int_0^{c_i} f(x_0 + th, y(x_0 + th))dt,$$

а интеграл приблизим квадратурной формулой *по тем же самым узлам* $\{c_1, \dots, c_s\}$, *но с другими коэффициентами* $\{a_{i1}, \dots, a_{is}\}$:

$$y(x_0 + c_i h) \approx y_0 + h \sum_{j=1}^s a_{ij} f(x_0 + c_j h, y(x_0 + c_j h)).$$

Обозначая $y(x_0 + c_i h) \approx Y_i$, в итоге получим s -стадийный метод Рунге – Кутты общего вида:

$$\begin{aligned} Y_i &= y_0 + h \sum_{j=1}^s a_{ij} f(x_0 + c_j h, Y_j), \quad i = \overline{1, s}, \\ y_1 &= y_0 + h \sum_{i=1}^s b_i f(x_0 + c_i h, Y_i). \end{aligned} \quad (7.17)$$

Запись метода РК в виде (7.17) называется *симметричной*. Более распространенной является другая, эквивалентная форма записи, которая получается заменой переменных:

$$\kappa_i = f(x_0 + c_i h, Y_i), \quad (7.18)$$

или

$$Y_i = y_0 + h \sum_{j=1}^s a_{ij} \kappa_j. \quad (7.19)$$

В результате имеем

$$\begin{aligned} \kappa_i &= f\left(x_0 + c_i h, y_0 + h \sum_{j=1}^s a_{ij} \kappa_j\right), \quad i = \overline{1, s}, \\ y_1 &= y_0 + h \sum_{i=1}^s b_i \kappa_i. \end{aligned} \quad (7.20)$$

Таким образом, s -стадийный метод РК определяется $s^2 + 2s$ параметрами $\{a_{ij}\}_{i,j=1}^s$, $\{b_1, \dots, b_s\}$ и $\{c_1, \dots, c_s\}$. Традиционно эти параметры размещают в виде следующей таблицы:

$$\begin{array}{c|c} c & A \\ \hline b^T & \end{array} = \begin{array}{c|c} c_1 & a_{11} \dots a_{1s} \\ \vdots & \dots\dots\dots \\ c_s & a_{s1} \dots a_{ss} \\ \hline & b_1 \dots b_s \end{array}. \quad (7.21)$$

Матрицу $A = (a_{ij})$ называют *матрицей Бутчера*⁸, а всю таблицу в формуле (7.21) — *таблицей Бутчера*.

⁸Джон Бутчер (John Butcher, университет Окленда) — новозеландский математик, внесший большой вклад в развитие и популяризацию методов РК.

▷₁ Постройте таблицу Бутчера для методов Эйлера, а также для явного и неявного метода средних прямоугольников.

Заметим, что первая формула в (7.20), как и ее симметричный аналог из (7.17), задает систему нелинейных уравнений относительно неизвестных $\{\kappa_i\}_{i=1}^s$ ($\{Y_i\}_{i=1}^s$ соответственно), т. е. методы РК являются в общем случае неявными. Однако если коэффициенты метода удовлетворяют условиям

$$a_{ij} = 0 \quad \forall i \leq j$$

и

$$c_1 = 0,$$

то такой метод РК очевидно будет явным.

7.2.3. Условия порядка методов Рунге – Кутты

Неизвестные коэффициенты методов типа Рунге – Кутты (7.21) традиционно выбираются таким образом, чтобы получившийся метод имел как можно более высокий порядок точности.

Найдем условия, которым должен удовлетворять произвольный метод РК второго порядка. Рассмотрим приближенное решение y_1 как функцию переменной h . Поскольку функция метода имеет вид $\varphi(x_0 + h) = y_1$, то

$$\left. \frac{d^k \varphi}{dx^k} \right|_{x=x_0} = \left. \frac{d^k y_1}{dh^k} \right|_{h=0}.$$

Тогда для того чтобы метод имел порядок 2, по определению необходимо, чтобы выполнялись условия

$$y_1|_{h=0} = y_0, \tag{7.22a}$$

$$y_1'|_{h=0} = y'(x_0) = f, \tag{7.22б}$$

$$y_1''|_{h=0} = y''(x_0) = f_x + f_y f. \tag{7.22в}$$

Условие (7.22a) очевидно выполняется всегда. Дифференцируя вторую формулу из (7.20) по h , имеем

$$y_1' = \sum_i b_i \kappa_i + h \sum_i b_i \kappa_i', \tag{7.23}$$

$$y_1'' = 2 \sum_i b_i \kappa_i' + h \sum_i b_i \kappa_i''.$$

Во всех суммах здесь и далее суммирование идет от 1 до s . Используя (7.19), вычислим

$$\begin{aligned}\kappa'_i &= \frac{d}{dh} f(x_0 + c_i h, Y_i) = c_i f'_x() + f'_y() Y'_i = \\ &= c_i f'_x() + f'_y() \left(\sum_j a_{ij} \kappa_j + h \sum_j a_{ij} \kappa'_j \right).\end{aligned}\quad (7.24)$$

Здесь $() = (x_0 + c_i h, Y_i)$. Используя (7.23), из (7.23) получим

$$\begin{aligned}y'_1|_{h=0} &= \sum_i b_i \kappa_i|_{h=0} = \sum_i b_i f, \\ y''_1|_{h=0} &= 2 \sum_i b_i \kappa'_i|_{h=0} = 2 \sum_i b_i \left[c_i f_x + f_y f \sum_j a_{ij} \right].\end{aligned}$$

Сопоставляя эти формулы с (7.226), (7.22в), получим следующие условия второго порядка для методов РК:

$$\begin{aligned}\sum_i b_i &= 1, \\ \sum_i b_i c_i &= \frac{1}{2}, \\ \sum_i b_i \sum_j a_{ij} &= \frac{1}{2}.\end{aligned}\quad (7.25)$$

▷₂ Выведите условия третьего порядка.

Замечание 7.3. В подавляющем большинстве случаев при построении методов РК априори полагают $c_i = \sum_j a_{ij}$, так что последние два условия в (7.25) эквивалентны.

Замечание 7.4. При выводе условий (7.25) мы рассматривали случай скалярного ОДУ. Следует понимать, что в общем случае условия порядка для систем ОДУ не будут совпадать со скалярными.

7.2.4. Явные методы Рунге – Кутты

Общий вид

Явные методы Рунге – Кутты можно смело назвать классическими, так как именно такого рода формулы строили в свое время Рунге и Кутта.

Определение 7.5. Метод Рунге – Кутты называется явным, если его таблица Бутчера имеет вид, представленный в табл. 7.1.

Таблица 7.1

0					
c_2	a_{21}				
c_3	a_{31}	a_{32}			
\vdots	$\dots\dots\dots$				
c_s	a_{s1}	a_{s2}	\dots	$a_{s,s-1}$	
<hr/>					
	b_1	b_2	\dots	b_{s-1}	b_s

Общая схема вычисления $y_1 \approx y(x_0 + h)$ по формулам (7.20) для явных методов РК имеет простой вид:

$$\begin{aligned}
 \kappa_1 &= f(x_0, y_0) \\
 \kappa_2 &= f(x_0 + c_2 h, y_0 + h a_{21} \kappa_1), \\
 \kappa_3 &= f(x_0 + c_3 h, y_0 + h(a_{31} \kappa_1 + a_{32} \kappa_2)), \\
 &\dots\dots\dots \\
 \kappa_s &= f(x_0 + c_s h, y_0 + h(a_{s1} \kappa_1 + \dots + a_{s,s-1} \kappa_{s-1})); \\
 y_1 &= y_0 + h \sum_{i=1}^s b_i \kappa_i.
 \end{aligned} \tag{7.26}$$

Явные методы — наиболее простые в реализации методы типа РК: вспомогательные значения κ_i вычисляются последовательно, каждое последующее из них явно выражается через уже найденные. Такой процесс очевидно невозможен, если в верхнем треугольнике матрицы A (включая диагональ) есть ненулевые элементы.

Примеры явных методов Рунге — Кутты

Второй порядок. Используя условия (7.25), легко получить общий вид явных двустадийных методов второго порядка. Поскольку $c_1 = a_{11} = a_{12} = a_{22} = 0$, для оставшихся коэффициентов метода имеем условия

$$\begin{aligned}
 b_1 + b_2 &= 1, \\
 b_2 c_2 &= b_2 a_{21} = \frac{1}{2},
 \end{aligned}$$

откуда получим следующий общий вид таблицы Бутчера:

$$\begin{array}{c|c} 0 & \\ \beta^{-1}/2 & \beta^{-1}/2 \\ \hline & 1 - \beta \quad \beta \end{array} \quad (7.27)$$

Здесь $b_2 = \beta$ — параметр, как правило, лежащий в полуинтервале $(0, 1]$.

Третий порядок. Приведем в табл. 7.2 без вывода некоторые методы порядка 3.

Таблица 7.2

0				0			
1/3	1/3			1/2	1/2		
2/3	0	2/3		1	0	1	
				1	0	0	1
	1/4	0	3/4		1/6	1/6	0 1/6

▷₃ Постройте какой-нибудь трехстадийный явный метод РК третьего порядка.

Четвертый порядок. Наиболее популярный «классический» метод четвертого порядка определяется табл 7.3.

Таблица 7.3

0				
1/2	1/2			
1/2	0	1/2		
1	0	0	1	
<hr/>				
	1/6	2/6	2/6	1/6

7.3. Выбор шага численного интегрирования ОДУ

Пусть нам необходимо приближенно решить (или, как еще говорят, численно проинтегрировать) задачу Коши:

$$y'(x) = f(x, y(x)), \quad y(x_0) = y_0,$$

с помощью некоторого одношагового метода Φ . Как и ранее, приближенное значение решения в точке $x_i > x_0$ обозначаем y_i :

$$y_i \approx y(x_i).$$

Наиболее простой вариант реализации численного интегрирования состоит в использовании равномерной сетки узлов,

$$x_i = x_0 + ih,$$

где h — некоторое фиксированное число, которое также называют шагом интегрирования. В таком случае процесс решения имеет вид

$$y_{i+1} = \Phi(x_i, y_i, x_i + h), \quad i = 0, 1, 2, \dots$$

Такой подход обладает теми же недостатками, что и применение составных квадратурных формул с постоянным шагом: как правило, на практике невозможно «угадать» нужную величину шага, которая обеспечивала бы требуемую точность и, с другой стороны, не приводила бы к чрезмерным вычислительным затратам.

Поэтому численное интегрирование задачи Коши для ОДУ обычно проводят с адаптивной⁹ длиной шага:

$$x_{i+1} = x_i + h_{i+1}, \quad i = 0, 1, 2, \dots$$

где длина шага h_{i+1} зависит от результата вычислений на предыдущем, i -м шаге.

7.3.1. Адаптивный выбор шага

Существует два широко применяемых способа адаптивного выбора шага. Все они основаны на оценке константы погрешности. Первый способ — это уже знакомый нам метод двойного пересчета, или *правило Рунге*.

⁹С длиной шага, которая автоматически приспособливается к поведению решения.

Правило Рунге

Правило Рунге для численного решения ОДУ практически полностью аналогично случаю приближенного вычисления определенного интеграла. Рассмотрим первый шаг процесса численного интегрирования методом Φ порядка p . Выберем какую-то величину начального шага h и сделаем сначала один «большой» шаг длиной $2h$:

$$\tilde{y}_2 = \Phi(x_0, y_0, x_0 + 2h),$$

а также два шага длиной h :

$$y_1 = \Phi(x_0, y_0, x_0 + h), \quad y_2 = \Phi(x_0 + h, y_1, x_0 + 2h).$$

Наша цель — оценить погрешность приближенного решения в точке $x_0 + 2h$. Для этого сравним погрешности

$$\tilde{e}_2 = y(x_0 + 2h) - \tilde{y}_2$$

и

$$e_2 = y(x_0 + 2h) - y_2.$$

По определению локальной погрешности (7.11) имеем

$$\tilde{e}_2 = y(x_0 + 2h) - \tilde{y}_2 = C_0(2h)^{p+1} + O(h^{p+2}). \quad (7.28)$$

Выражение погрешности e_2 имеет более сложный вид (рис. 7.2). Оно состоит из двух частей:

$$e_2 = y(x_0 + 2h) - y_2 = r_2 + \varepsilon_2.$$

Для записи r_2 (локальной погрешности второго шага) необходимо ввести в рассмотрение u — функцию, удовлетворяющую исходному дифференциальному уравнению

$$u'(x) = f(x, u(x)),$$

но с начальным условием $u(x_0 + h) = y_1$. Тогда по определению имеем

$$r_2 = u(x_0 + 2h) - y_2 = C_1 h^{p+1} + O(h^{p+2}), \quad (7.29)$$

а

$$\varepsilon_2 = y(x_0 + 2h) - u(x_0 + 2h).$$

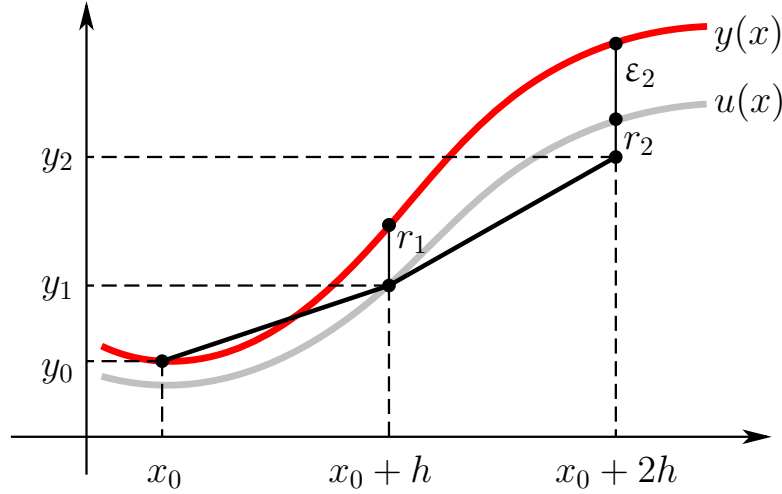


Рис. 7.2

Можно показать, что для величины ε_2 справедлива оценка

$$\varepsilon_2 = (1 + O(h))r_1,$$

где r_1 — локальная ошибка первого шага:

$$r_1 = y(x_0 + h) - y_1 = C_0 h^{p+1} + O(h^{p+2}),$$

а константа погрешности C_1 в (7.29) имеет вид

$$C_1 = C_0 + O(h).$$

Таким образом,

$$e_2 = \varepsilon_2 + r_2 = (r_1 + O(h^{p+2})) + ((C_0 + O(h))h^{p+1} + O(h^{p+2})),$$

откуда

$$e_2 = y(x_0 + 2h) - y_2 = 2C_0 h^{p+1} + O(h^{p+2}). \quad (7.30)$$

Вычитая из (7.28) выражение (7.30), получим

$$y_2 - \tilde{y}_2 = 2C_0(2^p - 1)h^{p+1} + O(h^{p+2}).$$

Это равенство позволяет вычислить оценку главной части погрешности для y_2 (см. (7.30)), которую традиционно обозначают err :

$$y(x_0 + 2h) = y_2 + \text{err} + O(h^{p+2}), \quad (7.31)$$

где

$$\text{err} = \frac{y_2 - \tilde{y}_2}{2^p - 1}, \quad (7.32)$$

а также приближенное значение константы погрешности

$$C_0 = \frac{\text{err}}{2h^{p+1}} + O(h). \quad (7.33)$$

Две последние формулы позволяют нам, во-первых, оценить погрешность e_2 , во-вторых, уточнить приближенное решение y_2 , и, в-третьих, выбрать оптимальный шаг интегрирования для дальнейших вычислений. Остановимся на каждом из этих пунктов подробнее.

Величина err , вычисляемая по формуле (7.32), по построению характеризует точность приближенного решения $y_2 \approx y(x_0 + 2h)$:

$$e_2 = \text{err} + O(h^{p+2}).$$

Требуемую точность в дальнейшем будем обозначать tol . Если $|\text{err}| < \text{tol}$, то полученное приближение приемлемо, в противном случае необходимо повторить вычисления с меньшим шагом (см. ниже).

Далее, согласно (7.31), величина

$$\hat{y}_2 = y_2 + \text{err} \quad (7.34)$$

является приближением к $y(x_0 + 2h)$ порядка $p + 1$:

$$y(x_0 + 2h) - \hat{y}_2 = O(h^{p+2}).$$

Таким образом, мы фактически получили способ повышения порядка метода на единицу. Величину \hat{y}_2 иногда называют *экстраполированным приближением*.

▷₁ Пользуясь формулами (7.34), (7.32), постройте общую формулу, по которой из любого одношагового метода Φ порядка p можно получить метод $\hat{\Phi}$ порядка $p + 1$.

▷₂ Запишите методы, полученные таким образом из явного и неявного методов Эйлера.

После того как вычислена оценка погрешности err , возможны два варианта развития событий в зависимости от выполнения неравенства

$$|\text{err}| < \text{tol}. \quad (7.35)$$

Предположим, что неравенство не выполняется, т. е. приближенное решение y_2 не является достаточно точным. Стандартный способ повышения точности в этом случае таков: y_2 «отбрасывается» и вычисления

повторяются с меньшим шагом \hat{h} . Выбор величины \hat{h} упрощается тем, что нам известна оценка константы погрешности (7.33). Согласно (7.30), мы должны выбрать \hat{h} из условия

$$|2C_0\hat{h}^{p+1}| < \text{tol},$$

откуда сразу же получаем

$$\hat{h} < \delta h, \quad (7.36)$$

где

$$\delta = \left(\frac{\text{tol}}{|\text{err}|} \right)^{\frac{1}{p+1}}. \quad (7.37)$$

Если же неравенство (7.35) выполняется, мы принимаем приближение y_2 (или его уточненное значение \hat{y}_2) и продолжаем процесс численного интегрирования из точки $x_0 + 2h$ с новым значением шага. Выбирают этот шаг из следующих соображений. Предположим, что константа погрешности на следующем шаге \tilde{C}_0 будет незначительно отличаться от C_0 (это будет верно, если шаг достаточно мал, а функция f достаточно гладкая). Тогда для того чтобы на новом шаге выполнялась оценка (7.35), мы совершенно аналогично предыдущему случаю должны выбрать \hat{h} из условий (7.36), (7.37). Только теперь величина δ будет больше единицы, т. е. шаг численного интегрирования *увеличится*.

Таким образом, **алгоритм автоматического (адаптивного) выбора шага численного интегрирования** имеет следующий общий вид:

- 1) $x \leftarrow x_0, y \leftarrow y_0, h \leftarrow h_0, X \leftarrow x_0 + H$;
- 2) если $x = X$, то завершаем алгоритм;
- 3) вычисляем $y_2 \leftarrow \Phi(x + h, \Phi(x, y, x + h), x + 2h)$,
 $\tilde{y}_2 \leftarrow \Phi(x, y, x + 2h)$;
- 4) находим оценку погрешности err (7.32) и коэффициент δ (7.37);
- 5) вычисляем $\hat{h} \leftarrow \alpha\delta h$, где $\alpha < 1$ — «страховочный множитель». Как правило, α выбирают в пределах от 0,7 до 0,9;
- 6) если $\delta < 1$, то полагаем $h \leftarrow \hat{h}$ и возвращаемся к пункту 3;
- 7) если же $\delta \geq 1$, то принимаем шаг: запоминаем пару значений $(x + 2h, y_2)$. Вместо y_2 здесь можно взять \hat{y}_2 ;
- 8) полагаем $x \leftarrow x + 2h, y \leftarrow y_2 (\hat{y}_2), h \leftarrow \min\{\hat{h}, \frac{X - x}{2}\}$;
- 9) возвращаемся к пункту 2.

Замечание 7.5. В силу погрешностей округления условие окончания алгоритма в пункте 2 следует использовать с осторожностью.

Замечание 7.6. В случае системы ОДУ оценка погрешности егг будет векторной величиной. Поэтому в формулах вида (7.37) следует использовать какую-либо векторную норму. Кроме этого, для некоторых задач величина относительной погрешности может быть более информативной, поэтому используют также формулы типа

$$\text{err} = \frac{1}{2^p - 1} \|D^{-1}(y_2 - \tilde{y}_2)\|,$$

где $D = \text{diag}(\tilde{y}_2)$. Если используется абсолютная погрешность, то $D = I$ — единичная матрица.

Замечание 7.7. Общепринято накладывать ограничения на величину шага h : он не должен быть слишком мал. Если же такой случай возникает, программы численного интегрирования обычно выдают сообщение об ошибке и прекращают работу. Кроме этого, накладываются ограничения на скорость увеличения шага, т. е. на величину δ :

$$\delta = \min \left(\delta_{\max}, \left(\frac{\text{tol}}{|\text{err}|} \right)^{\frac{1}{p+1}} \right).$$

Оценка погрешности с помощью вложенных методов

Данный способ оценки главной части погрешности более прост, но менее универсален, чем правило Рунге. Он требует наличия двух методов: метода Φ порядка p и метода $\hat{\Phi}$ порядка $q > p$. Обозначим

$$\begin{aligned} y_1 &= \Phi(x_0, y_0, x_0 + h), \\ \hat{y}_1 &= \hat{\Phi}(x_0, y_0, x_0 + h). \end{aligned}$$

По определению имеем

$$\begin{aligned} y(x_0 + h) - y_1 &= Ch^{p+1} + O(h^{p+2}), \\ y(x_0 + h) - \hat{y}_1 &= \hat{C}h^{q+1} + O(h^{q+2}). \end{aligned}$$

Отсюда с учетом того, что $q > p$, получим

$$\hat{y}_1 - y_1 = Ch^{p+1} + O(h^{p+2}).$$

Таким образом, величина

$$\text{err} = \hat{y}_1 - y_1 \tag{7.38}$$

с точностью до $O(h^{p+2})$ равна главной части локальной погрешности $y(x_0 + h) - y_1$. Дальнейшие рассуждения полностью аналогичны правилу Рунге. В частности, правило выбора оптимального шага (7.36), (7.37), а также общая схема алгоритма адаптивного выбора шага остаются без изменений. Единственное отличие заключается в способе оценки погрешности: вместо (7.32) имеем (7.38).

В заключение приведем один из наиболее популярных вложенных методов Рунге — Кутты — метод Дормана — Принса порядка 5(4) (табл. 7.4).

Таблица 7.4

0							
$\frac{1}{5}$	$\frac{1}{5}$						
$\frac{3}{10}$	$\frac{3}{40}$	$\frac{9}{40}$					
$\frac{4}{5}$	$\frac{44}{45}$	$-\frac{56}{15}$	$\frac{32}{9}$				
$\frac{8}{9}$	$\frac{19372}{6561}$	$-\frac{25360}{2187}$	$\frac{64448}{6561}$	$-\frac{212}{729}$			
1	$\frac{9017}{3168}$	$-\frac{355}{33}$	$\frac{46732}{5247}$	$\frac{49}{176}$	$-\frac{5103}{18656}$		
1	$\frac{35}{384}$	0	$\frac{500}{1113}$	$\frac{125}{192}$	$-\frac{2187}{6784}$	$\frac{11}{84}$	
\hat{b}_i	$\frac{35}{384}$	0	$\frac{500}{1113}$	$\frac{125}{192}$	$-\frac{2187}{6784}$	$\frac{11}{84}$	0
b_i	$\frac{5179}{57600}$	0	$\frac{7571}{16695}$	$\frac{393}{640}$	$-\frac{92097}{339200}$	$\frac{187}{2100}$	$\frac{1}{40}$

Удобство таких методов состоит в том, что для получения оценки погрешности требуется минимальное количество дополнительных вы-

числений: первая строчка коэффициентов \hat{b}_i в табл. 7.4 определяет приближение \hat{y}_1 порядка 5:

$$\hat{y}_1 = y_0 + h \sum_i \hat{b}_i \kappa_i,$$

а вторая — y_1 порядка 4 для оценки погрешности, т. е.

$$\text{err} = h \sum_i (\hat{b}_i - b_i) \kappa_i.$$

7.4. Многошаговые методы

7.4.1. Введение

Как и ранее, рассмотрим проблему численного решения задачи Коши

$$y'(x) = f(x, y(x)), \quad y(x_0) = y_0.$$

Предположим, что кроме начального значения y_0 нам известно еще и значение $y_1 = y(x_1)$, $x_1 = x_0 + h$. Каким образом можно использовать эту дополнительную информацию при вычислении приближения в точке $x_2 = x_1 + h$?

Проинтегрируем дифференциальное уравнение от x_1 до x_2 :

$$y(x_2) = y(x_1) + \int_{x_1}^{x_2} f(x, y(x)) dx = y_1 + h \int_0^1 F(t) dt,$$

где

$$F(t) = f(x_1 + th, y(x_1 + th)).$$

По условию нам известно, что

$$F(0) = f(x_1, y_1) = f_1, \quad F(-1) = f(x_0, y_0) = f_0,$$

поэтому можно приблизить подынтегральную функцию F интерполяционным многочленом P в форме Лагранжа:

$$P(t) = f_0 \Lambda_0(t) + f_1 \Lambda_1(t),$$

и получить приближенное равенство вида

$$y(x_2) \approx y_2 = y_1 + h \int_0^1 (f_0 \Lambda_0(t) + f_1 \Lambda_1(t)) dt = y_1 + h(\beta_0 f_0 + \beta_1 f_1),$$

где

$$\beta_0 = \int_0^1 \Lambda_0(t) dt = \int_0^1 \frac{t - 0}{-1 - 0} dt = -\frac{1}{2}, \quad \beta_1 = \int_0^1 \Lambda_1(t) dt = \int_0^1 \frac{t + 1}{0 + 1} dt = \frac{3}{2}.$$

В итоге получим формулу

$$y_2 = y_1 + \frac{h}{2}(3f_1 - f_0). \quad (7.39)$$

Это — *двухшаговый метод*, а точнее, двухшаговый явный метод Адамса. Использованное здесь обозначение

$$f_i = f(x_i, y_i),$$

будет активно применяться в дальнейшем.

Определение 7.6. Многошаговым (k -шаговым) методом решения задачи Коши называется отображение

$$\Phi : \left\{ \begin{bmatrix} x_0 \\ y_0 \end{bmatrix}, \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}, \dots, \begin{bmatrix} x_{k-1} \\ y_{k-1} \end{bmatrix}, x_k \right\} \mapsto y_k \approx y(x_k).$$

При $k = 1$ очевидно данное определение превращается в определение одношагового метода.

В общем случае задачи Коши нам известно только y_0 , следовательно, многошаговые методы при численной реализации нуждаются в процедуре вычисления «стартовых» значений y_1, \dots, y_{k-1} . Обычно для этого используются одношаговые методы. После этого уже возможно пошаговое применение метода Φ :

$$y_{n+1} = \Phi \left(\begin{bmatrix} x_{n-k+1} \\ y_{n-k+1} \end{bmatrix}, \dots, \begin{bmatrix} x_n \\ y_n \end{bmatrix}, x_{n+1} \right), \quad n = k-1, k, \dots \quad (7.40)$$

7.4.2. Явные методы Адамса

Нетрудно обобщить идею, которая использовалась при построении метода (7.39), на случай большего количества точек. Пусть известны значения y_0, y_1, \dots, y_{k-1} ,

$$y_i \approx y(x_i) = y(x_0 + ih).$$

Тогда имеем

$$y(x_k) = y(x_{k-1}) + h \int_0^1 F(t) dt, \quad (7.41)$$

$$F(t) = f(x_{k-1} + th, y(x_{k-1} + th)).$$

Функция F очевидно обладает свойством

$$F(0) = f_{k-1}, \quad F(-1) = f_{k-2}, \quad \dots, \quad F(1-k) = f_0,$$

или просто $F(t_j) = f_j$, где

$$t_j = 1 - k + j, \quad j = \overline{0, k-1}. \quad (7.42)$$

Строим для F интерполяционный многочлен P по узлам $\{t_j\}$:

$$F(t) \approx P(t) = \sum_{j=0}^{k-1} f_j \Lambda_j(t).$$

Заменяя F на P в интеграле (7.41), получим многошаговый метод вида

$$y_k = y_{k-1} + h \sum_{j=0}^{k-1} \beta_j f_j. \quad (7.43)$$

Коэффициенты β_j представляют собой интегралы от базисных функций Лагранжа Λ_i :

$$\beta_j = \int_0^1 \Lambda_j(t) dt = \int_0^1 \prod_{\substack{l=0 \\ l \neq j}}^{k-1} \frac{t - t_l}{t_j - t_l} dt, \quad (7.44)$$

т. е. по сути являются коэффициентами специальной интерполяционной квадратурной формулы, которая отличается от классических формул лишь тем, что использует узлы, лежащие за пределами отрезка интегрирования. Поэтому β_j можно найти не только по определению (7.44), но и используя тот факт, что по построению квадратурная формула

$$\int_0^1 F(t) dt \approx \sum_{j=0}^{k-1} \beta_j F(t_j)$$

должна иметь алгебраическую степень точности, равную $k - 1$. Подставляя в эту формулу $F(t) = t^m$, получим систему уравнений

$$\sum_{j=0}^{k-1} \beta_j (1 - k + j)^m = \frac{1}{m+1}, \quad m = 0, 1, \dots, k-1. \quad (7.45)$$

▷₁ Постройте явные методы Адамса для $k = 3$ и $k = 4$.

Замечание 7.8. В некоторых случаях строить формулы Адамса удобнее, применяя для представления интерполяционного многочлена P форму Ньютона.

7.4.3. Неявные методы Адамса

При интерполяции подынтегральной функции F в формуле (7.41) можно использовать не только известные значения F в точках $t = 0, -1, \dots, 1 - k$, но и неизвестное, «неявное» значение в точке $t = 1$, т. е.

$$F(1) = f(x_k, y_k).$$

Здесь y_k — искомое значение, которое нужно будет найти, решая полученное в итоге уравнение (систему уравнений). В результате такой модификации мы получим многошаговые методы Адамса, которые имеют общий вид

$$y_k = y_{k-1} + h \sum_{j=0}^k \hat{\beta}_j f_j. \quad (7.46)$$

Коэффициенты $\hat{\beta}_j$ вычисляются по формуле, аналогичной (7.44):

$$\hat{\beta}_j = \int_0^1 \hat{\Lambda}_j(t) dt = \int_0^1 \prod_{\substack{l=0 \\ l \neq j}}^k \frac{t - t_l}{t_j - t_l} dt. \quad (7.47)$$

Они могут быть найдены и как решение системы

$$\sum_{j=0}^k \hat{\beta}_j (1 - k + j)^m = \frac{1}{m + 1}, \quad m = 0, 1, \dots, k. \quad (7.48)$$

▷₂ Постройте двухшаговый и трехшаговый неявные методы Адамса.

7.4.4. Точность методов Адамса

Рассмотрим некоторый k -шаговый явный метод Адамса. Предположим, что все начальные значения y_0, y_1, \dots, y_{k-1} заданы точно, т. е.

$$y_j = y(x_j) \quad \forall j = \overline{0, k-1},$$

а y_k вычислено с помощью метода. Наша задача — оценить величину погрешности

$$y(x_k) - y_k.$$

По построению имеем

$$y(x_k) = y(x_{k-1}) + h \int_0^1 F(t) dt,$$

$$y_k = y_{k-1} + h \int_0^1 P(t) dt,$$

где

$$F(t) = f(x_{k-1} + th, y(x_{k-1} + th)),$$

а многочлен P интерполирует F по узлам $\{t_j\}_{j=0}^{k-1}$ (7.42). Отсюда получим

$$y(x_k) - y_k = h \int_0^1 r(t) dt,$$

где $r(t) = F(t) - P(t)$ — остаток полиномиальной интерполяции, который согласно формуле (5.28) можно представить в виде

$$r(t) = \frac{1}{k!} \frac{d^k}{dt^k} F(\xi) \underbrace{t(t+1) \dots (t+k-1)}_{\omega_k(t)},$$

где ξ — некоторая точка из интервала $(1-k, 0)$, зависящая от t . В силу знакопостоянства многочлена ω_k на отрезке $[0, 1]$ по уже неоднократно применявшейся нами теореме о среднем имеем

$$\int_0^1 r(t) dt = \frac{1}{k!} \frac{d^k}{dt^k} F(\eta) \int_0^1 t(t+1) \dots (t+k-1) dt = C \frac{d^k}{dt^k} F(\eta).$$

Здесь C — константа, зависящая только от k , а $\eta \in (1-k, 0)$. Дифференцируя F как сложную функцию, нетрудно убедиться в том, что

$$\frac{d^k}{dt^k} F(\eta) = O(h^k).$$

Отсюда окончательно получим

$$y(x_k) - y_k = hC \frac{d^k}{dt^k} F(\eta) = O(h^{k+1}),$$

т. е. явный k -шаговый метод Адамса имеет порядок точности, равный k .

В случае неявного k -шагового метода Адамса интерполяционный многочлен P имеет степень $k + 1$. Проводя полностью аналогичные рассуждения, легко показать, что неявный k -шаговый метод Адамса имеет порядок точности, равный $k + 1$.

7.5. Численное решение краевых задач. Метод стрельбы

7.5.1. Двухточечные краевые задачи

Знакомство с краевыми задачами для систем ОДУ начнем с рассмотрения примера. В плоскости xOy рассмотрим полет снаряда, выпущенного из точки $(0, 0)$ с начальной абсолютной скоростью v_0 под углом α к оси Ox . Траектория полета описывается системой ОДУ

$$\begin{cases} y'(x) = \operatorname{tg} \theta(x), \\ v'(x) = -g \frac{\operatorname{tg} \theta(x)}{v(x)} - k \frac{v(x)}{\cos \theta(x)}, \\ \theta'(x) = -\frac{g}{v(x)^2} \end{cases} \quad (7.49)$$

с начальными условиями

$$y(0) = 0, \quad v(0) = v_0, \quad \theta(0) = \alpha. \quad (7.50)$$

Здесь $(x, y(x))$ — координаты снаряда, $v(x)$ — скорость снаряда, $\theta(x)$ — угол между вектором скорости и осью Ox в точке $(x, y(x))$, g — гравитационная постоянная (смасштабированная должным образом), k — коэффициент сопротивления воздуха.

Решить данную задачу Коши можно с помощью любого из рассмотренных ранее численных методов. Однако с практической точки зрения более интересна другая задача: *определить траекторию полета снаряда, при которой он поражает наземную цель, находящуюся в точке $x = X$* . Соответствующее решение уравнения (7.49) вместо начальных условий (7.50) должно удовлетворять условиям

$$y(0) = 0, \quad y(X) = 0, \quad v(0) = v_0, \quad (7.51)$$

которые называются *краевыми*. Уравнения (7.49) вместе с условиями (7.51) представляют собой типичный пример *двухточечной краевой задачи*.

Определение 7.7. Рассмотрим систему ОДУ

$$y'(x) = f(x, y(x)), \quad x \in [a, b], \quad (7.52a)$$

где y и f — вектор-функции:

$$y(x) = \begin{bmatrix} y^1(x) \\ \vdots \\ y^n(x) \end{bmatrix}, \quad f(x, y(x)) = \begin{bmatrix} f^1(x, y(x)) \\ \vdots \\ f^n(x, y(x)) \end{bmatrix}.$$

Зададим n условий (связей), которым должно удовлетворять решение системы (7.52а):

$$G_i(y(a), y(b)) = 0, \quad i = \overline{1, n}. \quad (7.52б)$$

Уравнения (7.52) определяют *двухточечную краевую задачу* для ОДУ общего вида.

Замечание 7.9. Помимо (вместо) краев отрезка $[a, b]$ в условиях (7.52б) могут фигурировать и другие точки из этого отрезка. В этом случае получим более общую задачу, которую называют *многоточечной задачей*.

Следует понимать, что краевая задача может иметь несколько или ни одного решения даже в том случае, когда соответствующая начальная задача (задача Коши) однозначно разрешима.

7.5.2. Метод стрельбы

Приступим к рассмотрению первого метода решения граничных задач вида (7.52) — метода стрельбы. Это один из методов, основанных на сведении краевой задачи (которую мы не умеем решать) к последовательности задач Коши (которые мы решать уже научились).

Рассмотрим артиллерийскую краевую задачу (7.49), (7.51). Эту задачу можно переформулировать так: *найти такой угол наклона пушки $\alpha = \alpha^*$, при котором решение задачи Коши с условиями (7.50) проходит через точку $(X, 0)$* . Обозначим $y(x, \alpha)$ решение системы (7.49) с начальными условиями (7.50), зависящее от параметра α . Тогда алгоритм нахождения нужного угла α^* может выглядеть следующим образом:

- 1) выбираем два значения α_1 и α_2 , для которых заведомо выполняется условие $y(X, \alpha_1) < 0 < y(X, \alpha_2)$;
- 2) полагаем $\alpha \leftarrow (\alpha_1 + \alpha_2)/2$ и делаем пробный выстрел: решаем задачу Коши, находя $Y_\alpha = y(X, \alpha)$;
- 3) если $Y_\alpha = 0$ (или $|Y_\alpha| < \varepsilon$) — цель поражена, заканчиваем стрельбу;
- 4) если перелет ($Y_\alpha > 0$), полагаем $\alpha_2 \leftarrow \alpha$, в противном случае полагаем $\alpha_1 \leftarrow \alpha$;

5) переходим к шагу 2.

Это и есть простейший вариант метода стрельбы.

Описанная схема очевидно есть не что иное, как алгоритм метода бисекции для решения уравнения

$$F(\alpha) = y(X, \alpha) = 0.$$

Для ускорения сходимости можно применить более совершенные методы решения нелинейных уравнений, например метод секущих или Ньютона. Применение последнего, однако, затруднено необходимостью вычисления $F'(\alpha) = \frac{\partial}{\partial \alpha} y(X, \alpha)$.

Подведем локальный итог: *метод стрельбы заключается в подборе таких начальных условий, при которых решение данной системы ОДУ будет удовлетворять поставленным краевым условиям.*

Общая схема метода стрельбы

В общем случае ситуация осложняется тем, что неизвестных начальных условий может быть несколько. Например, вместо (7.51) можно рассмотреть условия

$$y(0) = 0, \quad y(X) = 0, \quad v(X) = v_X.$$

В этом случае нужно подбирать два начальных условия: $\theta(0)$ и $v(0)$, поэтому приходится использовать методы решения систем нелинейных уравнений: различные модификации метода Ньютона, метод Бroyдена и т. п. Запишем общую схему метода стрельбы в таком случае.

Пусть дана система ОДУ (7.52a), которую запишем в покомпонентной форме

$$\frac{d}{dx} y^i(x) = f^i(x, y^1(x), \dots, y^n(x)), \quad i = \overline{1, n}, \quad (7.53a)$$

и краевые условия, первые m из которых заданы в точке $x = b$, а остальные — в точке $x = a$, т. е. являются начальными:

$$\begin{cases} y^i(b) = \beta_i, & i = \overline{1, m}, \\ y^i(a) = \alpha_i, & i = \overline{m+1, n}. \end{cases} \quad (7.53б)$$

Для решения данной краевой задачи определим вектор-функцию

$$F = (F^1, \dots, F^m)^T : \mathbb{R}^m \rightarrow \mathbb{R}^m$$

следующим образом:

$$F^i(\alpha_1, \dots, \alpha_m) = y^i(b, \alpha_1, \dots, \alpha_m) - \beta_i, \quad i = \overline{1, m},$$

где $y^i(x, \alpha_1, \dots, \alpha_m)$ — решение исходной системы ОДУ (7.53a) с начальными условиями

$$y^i(a) = \alpha_i, \quad i = \overline{1, n} \quad (7.54)$$

(напомним, что начальные значения $\alpha_{m+1}, \dots, \alpha_n$ заданы по условию). Метод стрельбы решения краевой задачи (7.53) заключается в нахождении вектора $\alpha^* = (\alpha_1^*, \dots, \alpha_m^*)$, такого, что

$$F(\alpha^*) = 0.$$

Полученное решение $y(x, \alpha_1^*, \dots, \alpha_m^*)$ по построению будет являться решением краевой задачи (7.53).

Обратим особое внимание, что каждое вычисление $F(\alpha)$ требует решения задачи Коши (7.53a), (7.54), что весьма трудоемко. Более того, для реализации метода Ньютона, который считается наиболее эффективным методом решения систем уравнений, необходимо вычислять матрицу Якоби $\frac{\partial F}{\partial \alpha}$, вид которой, как правило, неизвестен. Эту матрицу приходится приближать с помощью конечных разностей, что, во-первых, ненадежно, и, во-вторых, очень накладно. Поэтому следует помнить, что метод стрельбы для краевых задач большой размерности очень трудоемок и ненадежен.

7.6. Решение краевых задач конечно-разностным методом

7.6.1. Общий нелинейный случай

Рассмотрим нелинейную краевую задачу для ОДУ второго порядка:

$$u''(x) = g(x, u(x), u'(x)), \quad u(a) = A, \quad u(b) = B. \quad (7.55)$$

Здесь $u : \mathbb{R} \rightarrow \mathbb{R}$, $g : \mathbb{R}^3 \rightarrow \mathbb{R}$. Метод конечных разностей или, как его чаще называют в русскоязычной литературе, *метод сеток* позволяет найти приближенное решение такой задачи в точках сетки $\{x_i\}_{i=0}^n$, которую чаще всего берут равномерной с шагом $h = (b - a)/n$:

$$x_i = a + ih, \quad i = \overline{0, n}. \quad (7.56)$$

Приближенное решение в этих узлах традиционно обозначают

$$y_i \approx u(x_i).$$

Поскольку по условию можно сразу положить $y_0 = A$, $y_n = B$, для нахождения $\{y_i\}_{i=1}^{n-1}$ необходимо из исходного дифференциального уравнения (7.55) получить $n - 1$ уравнений, которые бы связывали между собой эти неизвестные. Для этого необходимо выразить приближенные значения производных $u'(x_i)$ и $u''(x_i)$ через значения u в узлах сетки. Осуществить это можно разными способами, но все они в конечном итоге приводят к одинаковым результатам. Мы воспользуемся способом интерполяции.

Разностные производные

Для приближения производных применяется та же идея, что и при аппроксимации интегралов квадратурными формулами, а именно: *производная функции приближается производной многочлена, интерполирующего эту функцию по некоторому заранее заданному набору точек*.

Рассмотрим простейшие случаи. Пусть нам известны значения функции u в точках x_i и $x_{i+1} = x_i + h$. Интерполяционный многочлен в форме Ньютона по этим точкам имеет вид

$$P(x) = u(x_i) + \frac{u(x_{i+1}) - u(x_i)}{h}(x - x_i).$$

Отсюда получим

$$u'(x_i) \approx P'(x_i) = \frac{u(x_{i+1}) - u(x_i)}{h}. \quad (7.57)$$

Такая аппроксимация называется *правой разностной производной*.

Повторяя описанные действия для узлов интерполяции $\{x_{i-1}, x_i\}$, получим *левую разностную производную*

$$u'(x_i) \approx \frac{u(x_i) - u(x_{i-1}))}{h}. \quad (7.58)$$

Более точное приближение получается при использовании трех интерполяционных узлов $\{x_{i-1}, x_i, x_{i+1}\}$:

$$P(x) = u(x_{i-1}) + \frac{u(x_i) - u(x_{i-1}))}{h}(x - x_{i-1}) + u[x_{i-1}, x_i, x_{i+1}](x - x_{i-1})(x - x_i).$$

Дифференцируя этот многочлен, получим *центральную разностную производную*

$$u'(x_i) \approx \frac{u(x_{i+1}) - u(x_{i-1}))}{2h} \quad (7.59)$$

и *вторую разностную производную*

$$u''(x_i) \approx \frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))}{h^2}. \quad (7.60)$$

▷₁ Выведите эти формулы.

Разностная схема

Приближая производные в дифференциальном уравнении (7.55) по формулам (7.59) и (7.60), мы получим набор приближенных равенств вида

$$\frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))}{h^2} \approx g\left(x_i, u(x_i), \frac{u(x_{i+1}) - u(x_{i-1}))}{2h}\right).$$

Здесь индекс i изменяется от 1 до $n - 1$. Теперь заменим $u(x_i)$ на их приближенные значения y_i и окончательно получим следующую систему нелинейных уравнений, а точнее семейство систем, зависящих от параметра n :

$$\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} = g\left(x_i, y_i, \frac{y_{i+1} - y_{i-1}}{2h}\right), \quad i = \overline{1, n-1}, \quad (7.61a)$$

$$y_0 = A, \quad y_n = B. \quad (7.61b)$$

Такое семейство систем уравнений называется *разностной схемой*. Решая систему любым из известных нам численных методов при заданном значении шага h (числа n), получим приближенное решение краевой задачи (7.55).

7.6.2. Линейный случай

Наиболее просто конечно-разностный метод реализуется и исследуется в случае, когда уравнение линейно. Для простоты мы рассмотрим уравнение теплопроводности, которое описывает стационарное распределение тепла в стержне, на концах которого поддерживается постоянная температура A и B соответственно.

$$-u''(x) + q(x)u(x) = f(x), \quad u(a) = A, \quad u(b) = B. \quad (7.62)$$

Здесь $q(x)$ — коэффициент теплоотдачи, $f(x)$ — плотность источников тепла в точке x . Решение этого уравнения — функция u — определяет температуру стержня в точке x .

В данном случае уравнения (7.61a) превращаются в систему линейных уравнений

$$-\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} + q_i y_i = f_i, \quad (7.63)$$

где $f_i = f(x_i)$; $q_i = q(x_i)$. Ее можно переписать в виде

$$-y_{i-1} + (2 + q_i h^2)y_i - y_{i+1} = h^2 f_i. \quad (7.64)$$

Обозначая

$$d_i = 2 + q_i h^2,$$

а также учитывая краевые условия $y_0 = A$, $y_n = B$, окончательно получим СЛАУ

$$L_h y_h = f_h, \quad (7.65)$$

где

$$L_h = \frac{1}{h^2} \begin{bmatrix} h^2 & 0 & & & \\ -1 & d_1 & -1 & & \\ & -1 & d_2 & -1 & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & d_{n-2} & -1 \\ & & & & -1 & d_{n-1} & -1 \\ & & & & & 0 & h^2 \end{bmatrix}; \quad (7.66)$$

$y_h = (y_0, y_1, \dots, y_{n-1}, y_n)^T$; $f_h = (A, f_1, \dots, f_{n-1}, B)^T$. Нижний индекс h при записи системы (7.65) призван акцентировать факт зависимости размерности этой задачи и ее коэффициентов от величины шага сетки.

7.6.3. Устойчивость и сходимость метода сеток

Исследуем сходимость метода сеток для задачи (7.62). Будем считать $q(x) > 0 \forall x \in [a, b]$, так что имеем

$$d_i = 2 + q_i h^2 > 2 \quad \forall i = \overline{1, n-1}. \quad (7.67)$$

Это свойство гарантирует диагональное преобладание для матрицы L_h , а значит, существование и единственность решения СЛАУ (7.65), а также применимость метода прогонки для ее решения.

Теперь наша задача — доказать, что метод сеток сходится, т. е.

$$\max_{0 \leq i \leq n} |u(x_i) - y_i| = \|u_h - y_h\| \xrightarrow{h \rightarrow 0} 0. \quad (7.68)$$

Здесь и далее $u_h = (u(x_0), u(x_1), \dots, u(x_{n-1}), u(x_n))^T$ — вектор точных значений решения в узлах сетки, $\|\cdot\|$ — любая векторная норма.

Рассмотрим невязку, полученную при подстановке u_h вместо y_h в систему (7.65):

$$\psi_h = L_h u_h - f_h.$$

Такой вектор принято называть *погрешностью аппроксимации*. Поскольку по условию $L_h y_h = f_h$, имеем

$$L_h(u_h - y_h) = \psi_h,$$

откуда

$$\|u_h - y_h\| \leq \|L_h^{-1}\| \|\psi_h\|.$$

Из последнего соотношения становятся очевидными **достаточные условия сходимости метода**: для того чтобы выполнялось (7.68), достаточно выполнения следующих двух условий:

- 1) $\|\psi_h\| \xrightarrow{h \rightarrow 0} 0$ (условие аппроксимации);
- 2) $\|L_h^{-1}\| \leq M < \infty \forall h \leq h_0$, где M не зависит от h (условие устойчивости), а h_0 — некоторое положительное число.

Аппроксимация

Докажем выполнение первого условия сходимости. Рассмотрим ψ_i — i -ю компоненту вектора $\psi_h = L_h u_h - f_h$. Пусть $i = \overline{1, n-1}$, тогда

$$\begin{aligned}\psi_i &= -\frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} + q_i u_i - f_i = \\ &= -\frac{u(x_i - h) - 2u(x_i) + u(x_i + h)}{h^2} + q(x_i)u(x_i) - f(x_i) = \\ &= u''(x_i) - \frac{u(x_i - h) - 2u(x_i) + u(x_i + h)}{h^2}.\end{aligned}$$

Раскладывая в ряд Тейлора величины $u(x_i \pm h)$, в итоге получим отсюда

$$\psi_i = O(h^2), \quad i = \overline{1, n-1}.$$

▷₂ Прodelайте это.

Для $i = 0$ и $i = n$ имеем $\psi_i = 0$, так что окончательно

$$\|\psi_h\| = \max_i |\psi_i| = O(h^2) \xrightarrow{h \rightarrow 0} 0.$$

В таких случаях говорят, что *метод (схема) аппроксимирует дифференциальное уравнение со вторым порядком*.

Устойчивость

Для доказательства устойчивости, т. е. равномерной ограниченности нормы матрицы L_h^{-1} , достаточно показать, что если

$$L_h y_h = f_h,$$

то

$$\|y_h\| \leq M \|f_h\|, \quad M < \infty,$$

для всех h , меньших некоторого h_0 , причем константа M не должна зависеть от h . Можно показать, что в нашем случае

$$M = \left(1 + \frac{(b-a)^2}{8}\right)$$

(доказательство этого факта выходит за рамки нашего курса).

Таким образом, конечно-разностный метод для линейной краевой задачи теплопроводности (7.62) сходится.

Предметный указатель

- LDL^T -разложение, 35
- LU-разложение матрицы, 24
- QR-разложение матрицы, 44
- Арифметика с плавающей точкой, 9
- Барицентрическая формула Лагранжа, 138
- Главная часть остатка квадратур, 210
- Главный член локальной погрешности, 221
- Диагональное
 - преобладание, 28
 - — строгое, 28
- Задача
 - корректная, 13
 - некорректная, 13
- Задача интерполяции, 133
- Задача интерполяции кривой, 160
- Интерполяционный многочлен, 134
 - в барицентрической форме, 138
 - в форме Лагранжа, 136
 - в форме Ньютона, 139
- Квадратичная сходимость, 105
- Квадратурная сумма, 187
- Квадратурная формула, 187
 - Ньютона — Котеса, 197
 - Симпсона, 198
 - интерполяционная, 188
 - левых прямоугольников, 186
 - правых прямоугольников, 186
 - средних прямоугольников, 186
 - — составная, 207
 - трапеций, 187
 - — составная, 207
- Константа погрешности одношагового метода, 221
- Кратность собственного значения
 - алгебраическая, 76
 - геометрическая, 76
- Линейная сходимость, 105
- Мантисса, 6
- Матрица
 - диагонализируемая, 77
 - положительно определенная, 38
 - разреженная, 50
 - элементарного вращения, 93
- Матрица Вандермонда, 134
- Машинное число, 9
- Машинный эпсилон, 10
- Метод
 - Рунге — Кутты
 - — общего вида, 225
 - Гаусса, 21
 - Гаусса — Зейделя
 - — для СЛАУ, 57
 - — нелинейный, 119
 - Данилевского, 86
 - Якоби, 56
 - прогонки, 27
 - релаксации, 57
 - степенной метод, 77
- Многочлен
 - Бернштейна, 163
 - Лежандра, 202
 - наименее отклоняющийся от нуля, 146
- Норма
 - векторная
 - — p -норма, 15
 - — евклидова, 15
 - — максимум-норма, 15
 - матричная
 - — Фробениуса, 98

- индуцированная, 16
- подчиненная, 16
- спектральная, 16
- строчная максимум-норма, 16
- Обратная подстановка, 21
- Одношаговый метод, 218
- неявный, 219
- явный, 219
- Оператор
- простой структуры, 77
- Ортогональное преобразование, 40
- Ортогональные функции, 200
- Остаток интерполирования, 144
- Параметрические узлы, 160
- естественные, 161
- Плохо обусловленная задача, 14
- Погрешность
- аппроксимации, 252
- локальная, 220
- округления
- абсолютная, 10
- относительная, 10
- Погрешность интерполирования, 144
- Порядок сходимости, 105
- Правило
- округления, 7
- Проблема собственных значений, 77
- Прямой ход метода Гаусса, 21
- Разложение Холецкого, 39
- Разностная производная
- вторая, 250
- левая, 250
- правая, 250
- центральная, 250
- Разностная схема, 251
- Сверхлинейная сходимость, 105
- Следствие
- достаточное условие сходимости, 59
- Собственное значение, 58
- Собственное подпространство, 76
- Собственный вектор, 58
- Составная квадратурная формула, 206
- Спектр матрицы, 58
- Спектральный радиус, 58
- Сплайн
- бикубический, 180
- билинейный, 179
- Среднеквадратичное приближение
- дискретное, 168
- функции, 172
- Стандарт IEEE 754, 11
- Стационарный итерационный процесс, 55
- Теорема
- LDL^T -разложение, 35
- LU -разложение с выбором главного элемента, 33
- QR -алгоритм, упрощенная формулировка, 90
- Достаточное условие локального минимума, 128
- К. Ф. Гаусс, 201
- Лагранжа о среднем значении, 123
- Ролля, 145
- критерий сходимости стационарных итерационных процессов, 58
- о повышенной АСТ симметричных КФ, 195
- о среднем, 190
- принцип сжимающих отображений, 120
- разложение Холецкого, 39
- связь метода Гаусса и LU -разложения, 24
- Точка интерполяции, 133
- Узел интерполяции, 133
- Условие Липшица, 120
- Условия Вольфе, 130
- Форма Фробениуса, 86
- Форма Хессенберга, 69
- Формат
- CSC, 51
- CSR, 51
- MSR, 52
- координатный, 50
- Фундаментальный базис, 134
- двумерной алгебраической интерполяции, 178
- Функция
- билинейная, 179
- интегрируемая с квадратом, 171
- Характеристический многочлен, 76
- Чебышевские узлы, 147
- Число

- обусловленности
- — матрицы, [17](#)
- с плавающей точкой, [6](#)
- — денормализованное, [9](#)
- — нормализованное, [6](#)

Экспонента числа с плавающей точкой, [6](#)

Элементарная матрица

- второго типа, [85](#)
- первого типа, [85](#)

Список рекомендуемой литературы

- [1] *Бахвалов, Н. С.* Численные методы в задачах и упражнениях : учеб. пособие / Н. С. Бахвалов, А. В. Лапин, Е. В. Чижонков. — М., 2000. — 190 с.
- [2] *Богачев, К. Ю.* Практикум на ЭВМ. Методы приближения функций : учеб. пособие / К. Ю. Богачев, 3-е изд. — М., 2002. — 190 с.
- [3] *Богачев, К. Ю.* Практикум на ЭВМ. Методы решения линейных систем и нахождения собственных значений : учеб. пособие / К. Ю. Богачев, 2-е изд. — М., 1999. — 200 с.
- [4] *Голуб, Дж.* Матричные вычисления : пер. с англ / Дж. Голуб, Ч. Ван Лоун. — М., 1999. — 548 с.
- [5] *Икрамов, Х. Д.* Несимметричная проблема собственных значений. Численные методы / Х. Д. Икрамов. — М., 1991. — 240 с.
- [6] Алгоритмы: Построение и анализ / Т. Кормен [и др.]; 2-е изд. М., 2012. — 1293 с.
- [7] *Крылов, В. И.* Приближенное вычисление интегралов / В. И. Крылов. — М., 1967. — 500 с.
- [8] *Крылов, В. И.* Вычислительные методы высшей математики : в 2 т./ В. И. Крылов, В. В. Бобков, П. И. Монастырный; Т. 1. Под ред. И. П. Мысовских. — Минск, 1972. — 584 с.
- [9] *Ланцош, К.* Практические методы прикладного анализа : пер. с англ. / К. Ланцош — М., 1961. — 524 с.
- [10] *Ляшко, И. И.* Методы вычислений / И. И. Ляшко, В. Л. Макаров, А. А. Скоробогатько — Киев, «Вища школа», 1977. — 408 с.

- [11] Самарский, А. А. Численные методы : учеб. пособие / А. А. Самарский, А. В. Гулин. — М., 1989. — 432 с.
- [12] Тыртышников, Е. Е. Матричный анализ и линейная алгебра : учеб. пособие / Е. Е. Тыртышников. — М., 2007. — 480 с.
- [13] Уилкинсон, Дж. Алгебраическая проблема собственных значений : пер. с англ. / Дж. Уилкинсон. — М., 1970. — 564 с.
- [14] Фаддеев, Д. К. Вычислительные методы линейной алгебры / Д. К. Фаддеев, В. К. Фаддеева. — М., 1963. — 658 с.
- [15] Цехан, О. Б. Матричный анализ : учеб. пособие / О. Б. Цехан. — Гродно, 2010. — 371 с.
- [16] Berrut, J.-P. Barycentric lagrange interpolation / J.-P. Berrut, L. N. Trefethen // SIAM Rev. 2004. №46. P. 501–517.
- [17] Goldberg, D. “What every computer scientist should know about floating-point arithmetic” / D. Goldberg // ACM Computing Surveys. 1991. Vol. 23 №1. P. 5–48.
- [18] Higham, N. J. Accuracy and stability of numerical algorithms / N. J. Higham, 2nd ed. Philadelphia, 2002. — 680 с.
- [19] Numerical Recipes in C: the art of scientific computing / W. H. Press [et al.]. — 2nd ed. Cambridge, 1997. — 994 p.
- [20] Saad, Y. Iterative methods for sparse linear systems / Y. Saad- 2nd ed. Philadelphia, 2003. — 528 p.
- [21] Sewell, G. Computational Methods of Linear Algebra / G. Sewell. — 2nd ed. New Jersey, 2005. — 268 p.
- [22] Trefethen, L. N. Numerical Linear Algebra / L. N. Trefethen, D. Bau. Philadelphia, 1997. — 373 p.

