

С.П. Шарый

Курс
ВЫЧИСЛИТЕЛЬНЫХ
МЕТОДОВ

Курс

ВЫЧИСЛИТЕЛЬНЫХ
МЕТОДОВ

С. П. ШАРЫЙ

Институт вычислительных технологий СО РАН

Новосибирск – 2015

Книга является систематическим учебником по курсу вычислительных методов и написана на основе лекций, читаемых автором на механико-математическом факультете Новосибирского государственного университета. Особенностью книги является изложение методов интервального анализа и результатов конструктивной математики, связанных с традиционными разделами численного анализа.

Оглавление

Предисловие	8
Глава 1. Введение	9
1.1 Погрешности вычислений	11
1.2 Компьютерная арифметика	16
1.3 Обусловленность математических задач	18
1.4 Интервальная арифметика	21
1.5 Интервальные расширения функций	27
1.6 Элементы конструктивной математики	31
1.7 Сложность задач и трудоёмкость алгоритмов	33
1.8 Доказательные вычисления на ЭВМ	35
Литература к главе 1	38
Глава 2. Численные методы анализа	41
2.1 Введение	41
2.2 Интерполирование функций	44
2.2а Постановка задачи и её свойства	44
2.2б Интерполяционный полином Лагранжа	49
2.2в Разделённые разности и их свойства	52
2.2г Интерполяционный полином Ньютона	59
2.2д Погрешность алгебраической интерполяции	63
2.3 Полиномы Чебышёва	68
2.3а Определение и основные свойства	68
2.3б Применения полиномов Чебышёва	73
2.4 Интерполяция с кратными узлами	75
2.5 Общие факты алгебраической интерполяции	80
2.6 Сплайны	86

2.6a	Элементы теории	86
2.6б	Интерполяционные кубические сплайны	89
2.6в	Экстремальное свойство кубических сплайнов	96
2.7	Нелинейные методы интерполяции	97
2.8	Численное дифференцирование	99
2.8a	Интерполяционный подход	100
2.8б	Оценка погрешности дифференцирования	105
2.8в	Метод неопределённых коэффициентов	112
2.8г	Полная погрешность дифференцирования	114
2.9	Алгоритмическое дифференцирование	118
2.10	Приближение функций	120
2.10a	Обсуждение постановки задачи	120
2.10б	Существование и единственность решения задачи приближения	123
2.10в	Задача приближения в евклидовом пространстве	126
2.10г	Среднеквадратичное приближение функций	129
2.11	Полиномы Лежандра	134
2.11a	Мотивация и определение	134
2.11б	Основные свойства полиномов Лежандра	139
2.12	Численное интегрирование	143
2.12a	Постановка и обсуждение задачи	143
2.12б	Простейшие квадратурные формулы	147
2.12в	Квадратурная формула Симпсона	151
2.12г	Интерполяционные квадратурные формулы	158
2.12д	Дальнейшие формулы Ньютона-Котеса	161
2.12е	Метод неопределённых коэффициентов	165
2.13	Квадратурные формулы Гаусса	166
2.13a	Задача оптимизации квадратур	166
2.13б	Простейшие квадратуры Гаусса	168
2.13в	Выбор узлов для квадратурных формул Гаусса	172
2.13г	Практическое применение формул Гаусса	175
2.13д	Погрешность квадратур Гаусса	178
2.14	Составные квадратурные формулы	181
2.15	Сходимость квадратур	184
2.16	Вычисление интегралов методом Монте-Карло	189
2.17	Правило Рунге для оценки погрешности	194
	Литература к главе 2	195

Глава 3. Численные методы линейной алгебры	199
3.1 Задачи вычислительной линейной алгебры	199
3.2 Теоретическое введение	202
3.2а Основные понятия	202
3.2б Собственные числа и собственные векторы	209
3.2в Разложения матриц, использующие спектр	213
3.2г Сингулярные числа и сингулярные векторы	215
3.2д Сингулярное разложение матриц	220
3.2е Матрицы с диагональным преобладанием	223
3.3 Нормы векторов и матриц	225
3.3а Векторные нормы	225
3.3б Топология на векторных пространствах	230
3.3в Матричные нормы	235
3.3г Подчинённые матричные нормы	239
3.3д Топология на множествах матриц	244
3.3е Энергетическая норма	246
3.3ж Спектральный радиус	248
3.3з Матричный ряд Неймана	253
3.4 Приложения сингулярного разложения	256
3.4а Исследование неособенности и ранга матриц	256
3.4б Решение систем линейных уравнений	257
3.4в Малоранговые приближения матрицы	258
3.4г Метод главных компонент	260
3.5 Обусловленность систем линейных уравнений	262
3.5а Число обусловленности матриц	262
3.5б Примеры плохообусловленных матриц	267
3.5в Практическое применение числа обусловленности	269
3.6 Прямые методы решения линейных систем	273
3.6а Решение треугольных линейных систем	276
3.6б Метод Гаусса для решения линейных систем	277
3.6в Матричная интерпретация метода Гаусса	280
3.6г Метод Гаусса с выбором ведущего элемента	283
3.6д Существование LU-разложения	287
3.6е Разложение Холецкого	291
3.6ж Метод Холецкого	293
3.7 Методы на основе ортогональных преобразований	298
3.7а Обусловленность и матричные преобразования	298
3.7б QR-разложение матриц	301
3.7в Ортогональные матрицы отражения	303

3.7г	Метод Хаусхолдера	307
3.7д	Матрицы вращения	311
3.7е	Процессы ортогонализации	314
3.8	Метод прогонки	318
3.9	Стационарные итерационные методы	324
3.9а	Краткая теория	324
3.9б	Сходимость стационарных одношаговых методов	327
3.9в	Подготовка системы к итерационному процессу	334
3.9г	Оптимизация скалярного предобуславливателя	337
3.9д	Итерационный метод Якоби	340
3.9е	Итерационный метод Гаусса-Зейделя	345
3.9ж	Методы релаксации	350
3.10	Нестационарные итерационные методы	355
3.10а	Теоретическое введение	355
3.10б	Метод наискорейшего спуска	360
3.10в	Метод минимальных невязок	366
3.10г	Метод сопряжённых градиентов	369
3.11	Методы установления	373
3.12	Теория А.А. Самарского	375
3.13	Вычисление определителей и обратных матриц	379
3.14	Оценка погрешности приближённого решения	382
3.15	Линейная задача о наименьших квадратах	385
3.16	Проблема собственных значений	386
3.16а	Обсуждение постановки задачи	386
3.16б	Обусловленность проблемы собственных значений	389
3.16в	Коэффициенты перекоса матрицы	393
3.16г	Круги Гершгорина	397
3.16д	Отношение Рэлея	400
3.17	Численные методы для проблемы собственных значений	403
3.17а	Предварительное упрощение матрицы	403
3.17б	Степенной метод	405
3.17в	Обратные степенные итерации	413
3.17г	Сдвиги спектра	415
3.17д	Метод Якоби	417
3.17е	Базовый QR-алгоритм	423
3.17ж	Модификации QR-алгоритма	426
3.18	Численные методы сингулярного разложения	428
	Литература к главе 3	429

Глава 4. Решение нелинейных уравнений и их систем	434
4.1 Введение	434
4.2 Вычислительно-корректные задачи	436
4.2а Предварительные сведения и определения	436
4.2б Задача решения уравнений не является вычислительно-корректной	439
4.2в ε -решения уравнений	441
4.2г Недостаточность ε -решений	443
4.3 Векторные поля и их вращение	446
4.3а Векторные поля	446
4.3б Вращение векторных полей	448
4.3в Индексы особых точек	451
4.3г Устойчивость особых точек	452
4.3д Вычислительно-корректная постановка	454
4.4 Классические методы решения уравнений	455
4.4а Предварительная локализация решений	456
4.4б Метод дихотомии	457
4.4в Метод простой итерации	460
4.4г Метод Ньютона и его модификации	463
4.4д Методы Чебышёва	465
4.5 Классические методы решения систем уравнений	467
4.5а Метод простой итерации	467
4.5б Метод Ньютона и его модификации	468
4.6 Интервальные линейные системы уравнений	470
4.7 Интервальные методы решения уравнений	472
4.7а Основы интервальной техники	472
4.7б Одномерный интервальный метод Ньютона	475
4.7в Многомерный интервальный метод Ньютона	478
4.7г Метод Кравчика	481
4.8 Глобальное решение уравнений и систем	482
Литература к главе 4	487
Обозначения	491
Краткий биографический словарь	495
Предметный указатель	502

Предисловие

Представляемая вниманию читателей книга написана на основе курса лекций по вычислительным методам, которые читаются автором на механико-математическом факультете Новосибирского государственного университета. Её содержание в основной своей части традиционно и повторяет на современном уровне тематику, заданную ещё в знаменитых «Лекциях о приближённых вычислениях» акад. А.Н. Крылова, первом в мире систематическом учебнике методов вычислений. Условно материал книги можно назвать «вычислительные методы-1», поскольку в стандарте университетского образования, существует вторая часть курса, посвящённая численному решению дифференциальных уравнений, как обыкновенных, так и в частных производных, интегральных уравнений и др.

Вместе с тем, книга имеет ряд особенностей. Во-первых, в ней широко представлены элементы интервального анализа и современные интервальные методы для решения традиционных задач вычислительной математики. Во-вторых, автор счёл уместным поместить в книгу краткий очерк идей конструктивной математики и теории сложности вычислений, тесно связанных с предметом математики вычислительной.

Глава 1

Введение

Курс методов вычислений является частью более широкой математической дисциплины — вычислительной математики, которую можно неформально определить «как математику вычислений» или «математику, возникающую в связи с разнообразными процессами вычислений». При этом под «вычислениями» понимается не только получение числового ответа к задаче, доведение результата «до числа», но и получение конструктивных представлений (приближений) для различных математических объектов. С 70-х годов XX века, когда качественно нового уровня достигло развитие вычислительных машин и их применение во всех сферах жизни общества, можно встретить расширительное толкование содержания вычислительной математики, как «раздела математики, включающего круг вопросов, связанных с использованием ЭВМ» (определение А.Н. Тихонова).

Иногда в связи с вычислительной математикой и методами вычислений используют термин «численный анализ», возникший в США в конце 40-х годов XX века. Он более узок по содержанию, так как во главу угла ставит расчёты числового характера, а аналитические или символичные вычисления, без которых в настоящее время невозможно представить вычислительную математику и её приложения, отодвигает на второй план.

Развитие вычислительной математики в различные исторические периоды имело свои особенности и акценты. Начиная с античности (вспомним Архимеда) и вплоть до Нового времени вычислительные методы гармонично входили в сферу научных интересов крупнейших математиков — И. Ньютона, Л. Эйлера, Н.И. Лобачевского, К.Ф. Гаусса,

К.Г. Якоби и многих других, чьи имена остались в названиях популярных численных методов. В XX веке, и особенно в его второй половине, на первый план выдвинулась разработка и применение конкретных практических алгоритмов для решения сложных задач математического моделирования (в основном, вычислительной физики, механики и управления). Нужно было запускать и наводить ракеты, улучшать характеристики самолётов и других сложных технических устройств и т. п.

На развитие вычислительной математики большое или даже огромное влияние оказывали конкретные способы вычислений и вычислительные устройства, которые возникали по ходу развития технологий и применялись в процессах вычислений. В частности, огромное по своим последствиям влияние было испытано вычислительной математикой в середине XX века в связи с появлением электронных цифровых вычислительных машин.

Три типа задач, в основном, интересуют нас в связи с процессом вычислений:

- Как конструктивно найти (вычислить) тот или иной математический объект или его конструктивное приближение? К примеру, как найти производную, интеграл, решение дифференциального уравнения и т. п. вещи?
- Какова трудоёмкость нахождения тех или иных объектов? может ли она быть уменьшена и как именно?
- Если алгоритм для нахождения некоторого объекта уже известен, то как наилучшим образом организовать вычисления по этому алгоритму на том или ином конкретном вычислительном устройстве? Например, чтобы при этом уменьшить погрешность вычисления и/или сделать его менее трудоёмким?

Вопросы из последнего пункта сделались особенно актуальными в связи с развитием различных архитектур электронных вычислительных машин, в частности, с входением в нашу повседневную жизнь многопроцессорных и параллельных компьютеров.

Ясно, что все три отмеченных выше типа вопросов тесно связаны между собой. К примеру, если нам удаётся построить алгоритм для решения какой-либо задачи, то, оценив сложность его исполнения, мы тем самым предьявляем и верхнюю оценку трудоёмкости решения этой задачи.

Исторически сложилось, что исследования по второму пункту относятся, главным образом, к различным теориям вычислительной сложности и к теории алгоритмов, которая в 30-е годы XX века вычленилась из абстрактной математической логики. Но традиционная вычислительная математика, предметом которой считается построение и исследование конкретных численных методов, также немало способствует прогрессу в этой области.

Опять же, исторические и организационные причины привели к тому, что различные вычислительные методы для решения тех или иных конкретных задач относятся к другим математическим дисциплинам. Например, численные методы для отыскания экстремумов различных функций являются предметом вычислительной оптимизации, теории принятия решений и исследования операций.

1.1 Погрешности вычислений

Общеизвестно, что в практических задачах числовые данные почти всегда не вполне точны и содержат ошибки. Если эти данные являются, к примеру, результатами измерений, то за редким исключением они не могут быть произведены абсолютно точно.

Ошибкой или *погрешностью* приближённого значения \tilde{x} какой-либо величины называют разность между \tilde{x} и истинным значением x этой величины, т. е. $\tilde{x} - x$. Часто более удобно оперировать *абсолютной погрешностью* $\tilde{\Delta}$ приближённой величины, которая определяется как абсолютная величина погрешности, т. е.

$$\tilde{\Delta} = |\tilde{x} - x|, \quad (1.1)$$

поскольку во многих случаях знак погрешности неизвестен.

Практически точное значение интересующей нас величины x неизвестно, так что вместо точного значения абсолютной погрешности также приходится довольствоваться её приближёнными значениями. Оценку сверху для абсолютной погрешности называют *предельной* (или *граничной*) *абсолютной погрешностью*. В самом этом термине содержится желание иметь эту величину как можно более точной, т. е. как можно меньшей.

Таким образом, если Δ — предельная абсолютная погрешность значения \tilde{x} точной величины x , то

$$\tilde{\Delta} = |\tilde{x} - x| \leq \Delta,$$

и потому

$$\tilde{x} - \Delta \leq x \leq \tilde{x} + \Delta.$$

Вместо этого двустороннего неравенства удобно пользоваться следующей краткой и выразительной записью:

$$x = \tilde{x} \pm \Delta.$$

Фактически, вместо точного числа мы имеем здесь целый диапазон значений — числовой интервал $[\tilde{x} - \Delta, \tilde{x} + \Delta]$ возможных представителей для точного значения x .

Как правило, указание одной только абсолютной погрешности недостаточно для характеристики качества рассматриваемого приближения. Более полное понятие о нём можно получить из *относительной погрешности* приближения, которая определяется как отношение абсолютной погрешности к самому значению этой величины:

$$\tilde{\delta} = \frac{\tilde{\Delta}}{|x|}. \quad (1.2)$$

Относительная погрешность — безразмерная величина.

Предельной относительной погрешностью приближённого значения x называют число δ , оценивающее сверху его относительную погрешность. Как правило, и для абсолютной, и для относительной погрешностей в речи опускают эпитеты «предельная», поскольку именно предельные (граничные) погрешности являются реальными доступными нам (наблюдаемыми) величинами.

При записи приближённых чисел имеет смысл изображать их так, чтобы сама форма их написания давала характеристику об их точности. Ясно, что ненадёжные знаки представления чисел указывать смысла нет. Обычно принимают за правило писать числа так, чтобы все их значащие цифры кроме, может быть, последней были верны, а последняя цифра была бы сомнительной не более чем на единицу. Согласно этому правилу число 12340000, у которого цифра 4 уже сомнительна, нужно записывать в виде $1.23 \cdot 10^8$.

Значащей цифрой приближённого числа называется цифра в его представлении в заданной системе счисления, отличная от нуля, либо нуль, если он стоит между значащими цифрами или является представителем сохранённого разряда этого числа. Содержательное определение может состоять в том, что значащая цифра — это цифра из

представления числа, которая даёт существенную информацию о его относительной погрешности.

Значащие цифры могут быть верными или неверными.

Как изменяются абсолютные и относительные погрешности при выполнении арифметических операций с приближёнными числами? Приближённое число с заданной абсолютной погрешностью — это, фактически, целый интервал значений. По этой причине для абсолютных погрешностей поставленный вопрос решается формулами интервальной арифметики, рассматриваемой в §1.4. Здесь мы рассмотрим упрощённые версии этих операций.

Предложение 1.1.1 *Абсолютная погрешность суммы и разности приближённых чисел равна сумме абсолютных погрешностей операндов.*

Доказательство. Если x_1, x_2 — точные значения рассматриваемых чисел, \tilde{x}_1, \tilde{x}_2 — их приближённые значения, а Δ_1, Δ_2 — соответствующие предельные абсолютные погрешности, то

$$\tilde{x}_1 - \Delta_1 \leq x_1 \leq \tilde{x}_1 + \Delta_1, \quad (1.3)$$

$$\tilde{x}_2 - \Delta_2 \leq x_2 \leq \tilde{x}_2 + \Delta_2. \quad (1.4)$$

Складывая эти неравенства почленно, получим

$$(\tilde{x}_1 + \tilde{x}_2) - (\Delta_1 + \Delta_2) \leq x_1 + x_2 \leq (\tilde{x}_1 + \tilde{x}_2) + (\Delta_1 + \Delta_2).$$

Полученное соотношение означает, что величина $\Delta_1 + \Delta_2$ является предельной абсолютной погрешностью суммы $\tilde{x}_1 + \tilde{x}_2$.

Умножая обе части неравенства (1.4) на (-1) , получим

$$-\tilde{x}_2 - \Delta_2 \leq -x_2 \leq -\tilde{x}_2 + \Delta_2.$$

Складывая почленно с неравенством (1.3), получим

$$(\tilde{x}_1 - \tilde{x}_2) - (\Delta_1 + \Delta_2) \leq x_1 + x_2 \leq (\tilde{x}_1 - \tilde{x}_2) + (\Delta_1 + \Delta_2).$$

Отсюда видно, что величина $\Delta_1 + \Delta_2$ является предельной абсолютной погрешностью разности $\tilde{x}_1 - \tilde{x}_2$. ■

Для умножения и деления формулы преобразования абсолютной погрешности более громоздки. Точные результаты для операций между приближёнными величинами даются интервальной арифметикой, рассматриваемой ниже в §1.4.

Рассмотрим теперь эволюцию относительной погрешности в вычислениях.

Предложение 1.1.2 *Если слагаемые в сумме имеют одинаковый знак, то относительная погрешность суммы не превосходит наибольшей из относительных погрешностей слагаемых.*

Доказательство. Пусть складываются две приближённые величины, точные значения которых равны x_1 и x_2 , а относительные погрешности суть δ_1 и δ_2 . Тогда их абсолютные погрешности —

$$\Delta_1 = \delta_1 |x_1| \quad \text{и} \quad \Delta_2 = \delta_2 |x_2|.$$

Если $\delta = \max\{\delta_1, \delta_2\}$, то

$$\Delta_1 \leq \delta |x_1|, \quad \Delta_2 \leq \delta |x_2|.$$

Складывая полученные неравенства почленно, получим

$$\Delta_1 + \Delta_2 \leq \delta (|x_1| + |x_2|),$$

откуда

$$\frac{\Delta_1 + \Delta_2}{|x_1| + |x_2|} \leq \delta.$$

В числителе полученной дроби предельная абсолютная погрешность суммы, а в знаменателе — модуль точного значения суммы, если слагаемые имеют один и тот же знак. ■

Ситуация с относительной погрешностью принципиально меняется, когда в сумме слагаемые имеют разный знак, т. е. она является разностью. Если результат имеет меньшую абсолютную величину, чем абсолютные величины операндов, то значение дроби (1.2) возрастёт. А если вычитаемые числа очень близки друг к другу, то знаменатель в (1.2) сделается очень маленьким и относительная погрешность результата может катастрофически возрасти.

Пример 1.1.1 Рассмотрим вычитание чисел 1001 и 1000, каждое из которых является приближённым и известным с абсолютной точностью 0.1. Таким образом, относительные точности обоих чисел примерно равны 0.01%. Выполняя вычитание, получим результат 1, который имеет абсолютную погрешность $0.1 + 0.1 = 0.2$. Как следствие, относительная погрешность результата достигла 20%. ■

Предложение 1.1.3 Если погрешности приближённых чисел малы, то относительная погрешность их произведения приближённо (с точностью до членов более высокого порядка малости) равна сумме относительных погрешностей сомножителей.

Доказательство. Пусть x_1, x_2, \dots, x_n — точные значения рассматриваемых чисел, $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$ — их приближённые значения. Обозначим также $x := x_1 x_2 \dots x_n$, $\tilde{x} := \tilde{x}_1 \tilde{x}_2 \dots \tilde{x}_n$, и пусть $f(y_1, y_2, \dots, y_n) = y_1 y_2 \dots y_n$ — функция произведения n чисел. Разлагая её в точке (x_1, x_2, \dots, x_n) по формуле Тейлора с точностью до членов первого порядка, получим

$$\begin{aligned} \tilde{x} - x &= f(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n) - f(x_1, x_2, \dots, x_n) \\ &\approx \sum_{i=1}^n \frac{\partial f}{\partial y_i}(x_1, x_2, \dots, x_n) \cdot (\tilde{x}_i - x_i) \\ &= \sum_{i=1}^n x_1 \dots x_{i-1} x_{i+1} \dots x_n (\tilde{x}_i - x_i) \\ &= \sum_{i=1}^n x_1 x_2 \dots x_n \frac{\tilde{x}_i - x_i}{x_i}. \end{aligned}$$

Разделив на $x = x_1 x_2 \dots x_n$ обе части этого приближённого равенства и беря от них абсолютное значение, получим с точностью до членов второго порядка малости

$$\left| \frac{\tilde{x} - x}{x} \right| = \sum_{i=1}^n \left| \frac{\tilde{x}_i - x_i}{x_i} \right|,$$

что и требовалось. ■

Предложение 1.1.4 Если погрешности приближённых чисел малы, то относительная погрешность их частного приближённо (с точностью до членов более высокого порядка малости) равна сумме относительных погрешностей сомножителей.

Доказательство. Если $u = x/y$, то

$$\Delta u \approx \frac{\partial u}{\partial x} \Delta x + \frac{\partial u}{\partial y} \Delta y = \frac{\Delta x}{y} - \frac{x \Delta y}{y^2}.$$

Поэтому

$$\frac{\Delta u}{u} = \frac{\Delta x}{y \frac{x}{y}} - \frac{x \Delta y}{y^2 \frac{x}{y}} = \frac{\Delta x}{x} - \frac{\Delta y}{y},$$

так что

$$\left| \frac{\Delta u}{u} \right| \leq \left| \frac{\Delta x}{x} \right| + \left| \frac{\Delta y}{y} \right|.$$

Это и требовалось показать. ■

1.2 Компьютерная арифметика

Для правильного учёта погрешностей реализации вычислительных методов на различных устройствах и для правильной организации этих методов нужно знать детали конкретного способа вычислений. В современных электронных цифровых вычислительных машинах (ЭЦВМ), на которых выполняется подавляющая часть современных вычислений, эти детали реализации регламентируются специальным международным стандартом. Он был принят в 1985 году Институтом инженеров по электротехнике и электронике¹, профессиональной ассоциацией, объединяющей в своих рядах также специалистов по аппаратному обеспечению ЭВМ. Этот стандарт, коротко называемый IEEE 754, был дополнен и развит в 1995 году следующим стандартом IEEE 854 [26, 34], а затем в 2008 году появилась переработанная версия первого стандарта, которая получила наименование IEEE 754-2008.

Согласно этим стандартам вещественные числа представляются в ЭВМ в виде «чисел с плавающей точкой», в которых число хранится в форме мантиссы и показателя степени. Зафиксируем натуральное число β , которое будет называться основанием системы счисления. *Числами с плавающей точкой* называются числа вида

$$(\alpha_1 \beta^{-1} + \alpha_2 \beta^{-2} + \dots + \alpha_p \beta^{-p}) \cdot \beta^e,$$

которые условно можно записать в виде

$$0.\alpha_1 \alpha_2 \dots \alpha_p \cdot \beta^e,$$

где $0 \leq \alpha_i < \beta$, $i = 1, 2, \dots, p$. В выписанном представлении величина $0.\alpha_1 \alpha_2 \dots \alpha_p$ называется *мантиссой* числа, а p — количество значащих

¹Чаще всего его называют английской аббревиатурой IEEE от Institute of Electrical and Electronics Engineers.

цифр мантиссы — это *точность* рассматриваемой модели с плавающей точкой. На показатель степени e также обычно накладывается двустороннее ограничение $e_{\min} \leq e \leq e_{\max}$.

Стандарты IEEE 754/854 предписывают для цифровых ЭВМ значения $\beta = 2$ или $\beta = 10$, и в большинстве компьютеров используется $\beta = 2$, т. е. двоичная система счисления. С одной стороны, это вызвано особенностями физической реализации современных ЭВМ, где 0 соответствует отсутствию сигнала (заряда и т. п.), а 1 — его наличию. С другой стороны, двоичная система оказывается выгодной при выполнении с ней приближённых вычислений (см. [26]).

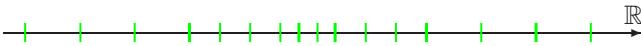


Рис. 1.1. Множество чисел, представимых в цифровой ЭВМ — дискретное конечное подмножество вещественной оси \mathbb{R} .

Как видим, числа с плавающей точкой обеспечивают практически фиксированную относительную погрешность представления вещественных чисел и изменяющуюся абсолютную погрешность.

Стандарты IEEE 754/854 предусматривают для чисел с плавающей точкой «одинарную точность» и «двойную точность», а также «расширенные» варианты этих представлений. При этом для хранения чисел одинарной точности отводится 4 байта памяти ЭВМ, для двойной точности — 8 байтов. Из этих 32 или 64 битов один бит зарезервирован для указания знака числа: 0 соответствует «−», а 1 соответствует «+». Таким образом, во внутреннем «машинном» представлении знак присутствует у любого числа, в том числе и у нуля.

Для двойной точности, наиболее широко распространённой в современных расчётах, диапазон чисел, представимых в ЭВМ простирается от примерно $2.22 \cdot 10^{-308}$ до $1.79 \cdot 10^{308}$. Помимо обычных чисел стандарты IEEE 754/854 описывают несколько специальных объектов вычислений. Это, прежде всего, машинная бесконечность и специальный нечисловой объект под названием NaN (названный как сокращение английской фразы «Not a Number»). NaN полезен во многих ситуациях, в частности, он может использоваться для сигнализации о нетипичных и исключительных событиях, случившихся в процессе вычислений, которые, тем не менее, нельзя было прерывать.

Очень важной характеристикой множества машинных чисел явля-

ется так называемое «машинное ε » (машинное эpsilon), которое характеризует густоту множества машинно-представимых чисел. Это наименьшее положительное число $\varepsilon_{\text{маш}}$, такое что в компьютерной арифметике $1 + \varepsilon_{\text{маш}} \neq 1$ при округлении к ближайшему. Из конструкции чисел с плавающей точкой следует тогда, что компьютер, грубо говоря, не будет различать чисел a и b , удовлетворяющих условию $1 < a/b < 1 + \varepsilon_{\text{маш}}$. Для двойной точности представления в стандарте IEEE 754/854 машинное эpsilon примерно равно $1.11 \cdot 10^{-16}$.

Принципиальной особенностью компьютерной арифметики, вызванной дискретностью множества машинных чисел и наличием округлений, является невыполнение некоторых общеизвестных свойств вещественной арифметики. Например, сложение чисел с плавающей точкой неассоциативно, т. е. в общем случае неверно, что

$$(a + b) + c = a + (b + c).$$

Читатель может проверить на любом компьютере, что в арифметике IEEE 754/854 двойной точности при округлении «к ближайшему»

$$(1 + 1.1 \cdot 10^{-16}) + 1.1 \cdot 10^{-16} \neq 1 + (1.1 \cdot 10^{-16} + 1.1 \cdot 10^{-16}).$$

Левая часть этого отношения равна 1, тогда как правая — ближайшему к единице справа машинно-представимому числу. Эта ситуация имеет место в любых приближённых вычислениях, которые сопровождаются округлениями, а не только при расчётах на современных цифровых ЭВМ.

Из отсутствия ассоциативности следует, что результат суммирования длинных сумм вида $x_1 + x_2 + \dots + x_n$ зависит от порядка, в котором выполняется попарное суммирование слагаемых, или, как говорят, от расстановки скобок в сумме. Каким образом следует организовывать такое суммирование в компьютерной арифметике, чтобы получать наиболее точные результаты? Ответ на этот вопрос существенно зависит от значений слагаемых, но в случае суммирования уменьшающихся по абсолютной величине величин суммировать нужно «с конца». Именно так, к примеру, лучше всего находить суммы большинства рядов.

1.3 Обусловленность математических задач

Вынесенный в заголовок этого параграфа термин — *обусловленность* — означает меру чувствительности решения задачи к изменениям (возмущениям) её входных данных. Ясно, что любая информация

подобного сорта чрезвычайно важна при практических вычислениях, так как позволяет оценивать достоверность результатов, полученных в условиях приближённого характера этих вычислений. С другой стороны, зная о высокой чувствительности решения мы можем предпринимать необходимые меры для компенсации этого явления — повышать разрядность вычислений, наконец, модифицировать или вообще сменить выбранный вычислительный алгоритм и т. п.

Существует несколько уровней рассмотрения поставленного вопроса. Во-первых, следует знать, является ли вообще непрерывной зависимость решения задачи от входных данных. Задачи, решение которых не зависит непрерывно от их данных, называют *некорректными*. Далее в §2.8г в качестве примера таких задач мы рассмотрим задачу численного дифференцирования. Во-вторых, в случае наличия этой непрерывности желательно иметь некоторую количественную меру чувствительности решения как функции от входных данных.

Переходя к формальным конструкциям, предположим, что в рассматриваемой задаче по значениям из множества \mathcal{D} входных данных мы должны вычислить решение задачи из множества ответов \mathcal{S} . Отображение $\phi : \mathcal{D} \rightarrow \mathcal{S}$, сопоставляющее всякому a из \mathcal{D} решение задачи из \mathcal{S} , мы будем называть *разрешающим отображением* (или *разрешающим оператором*). Отображение ϕ может быть выписано явным образом, если ответ к задаче задаётся каким-либо выражением. Часто разрешающее отображение задаётся неявно, как, например, при решении системы уравнений

$$F(a, x) = 0$$

с входными параметрами a . Даже при неявном задании нередко можно теоретически выписать вид разрешающего отображения, как, например, $x = A^{-1}b$ при решении системы линейных уравнений $Ax = b$ с квадратной матрицей A . Но в любом случае удобно предполагать существование этого отображения и некоторые его свойства. Пусть также \mathcal{D} и \mathcal{S} являются линейными нормированными пространствами. Очевидно, что самый первый вопрос, касающийся обусловленности задачи, требует, чтобы разрешающее отображение ϕ было непрерывным относительно некоторого задания норм в \mathcal{D} и \mathcal{S} .

Что касается числовой меры обусловленности математических задач, то существуют два подхода к её введению. Одни из них условно может быть назван *дифференциальным*, а другой основан на оценивании *константы Липшица* разрешающего оператора.

Пусть разрешающее отображение дифференцируемо по крайней мере в интересующей нас точке a из множества входных данных \mathcal{D} . Тогда можно считать, что

$$\phi(a + \Delta a) \approx \phi(a) + \phi'(a) \cdot \Delta a,$$

и потому мерой чувствительности решения может служить $\|\phi'(a)\|$. Для более детального описания зависимости различных компонент решения $\phi(a)$ от a часто привлекают отдельные частные производные $\frac{\partial \phi_i}{\partial a_j}$, т. е. элементы матрицы Якоби $\phi'(a)$ разрешающего отображения ϕ , которые при этом называют *коэффициентами чувствительности*. Интересна также мера относительной чувствительности решения, которую можно извлечь из соотношения

$$\frac{\phi(a + \Delta a) - \phi(a)}{\|\phi(a)\|} \approx \left(\frac{\phi'(a)}{\|\phi(a)\|} \cdot \|a\| \right) \frac{\Delta a}{\|a\|}.$$

Второй подход к определению обусловленности требует нахождения как можно более точных констант C_1 и C_2 в неравенствах

$$\|\phi(a + \Delta a) - \phi(a)\| \leq C_1 \|\Delta a\| \quad (1.5)$$

и

$$\frac{\|\phi(a + \Delta a) - \phi(a)\|}{\|\phi(a)\|} \leq C_2 \frac{\|\Delta a\|}{\|a\|}. \quad (1.6)$$

Величины этих констант, зависящие от задачи, а иногда и конкретных входных данных, берутся за меру обусловленности решения задачи.

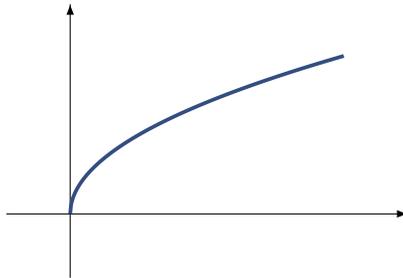


Рис. 1.2. Непрерывная функция $y = \sqrt{x}$ имеет бесконечную скорость роста при $x = 0$ и не является липшицевой

В связи с неравенствами (1.5)–(1.6) напомним, что вещественная функция $f : \mathbb{R}^n \supseteq D \rightarrow \mathbb{R}$ называется *непрерывной по Липшицу* (или просто *липшицевой*), если существует такая константа L , что

$$|f(x') - f(x'')| \leq L \cdot \text{dist}(x', x'') \quad (1.7)$$

для любых $x', x'' \in D$. Величину L называют при этом *константой Липшица* функции f на D . Понятие непрерывности по Липшицу формализует интуитивно понятное условие соразмерности изменения функции изменению аргумента. Именно, приращение функции не должно превосходить приращение аргумента (по абсолютной величине или в некоторой заданной метрике) более чем в определённое фиксированное число раз. При этом сама функция может быть и негладкой, как, например, модуль числа в окрестности нуля. Отметим, что понятие непрерывности по Липшицу является более сильным свойством, чем просто непрерывность или даже равномерная непрерывность, так как влечёт за собой их обоих.

Нетрудно видеть, что искомые константы C_1 и C_2 в неравенствах (1.5) и (1.6), характеризующие чувствительность решения задачи по отношению к возмущениям входных данных — это не что иное, как константы Липшица для разрешающего отображения ϕ и произведение константы Липшица L_ψ отображения $\psi : \mathcal{D} \rightarrow \mathcal{S}$, действующего по правилу $a \mapsto \phi(a)/\|\phi(a)\|$ на норму $\|a\|$. В последнем случае

$$\frac{\|\phi(a + \Delta a) - \phi(a)\|}{\|\phi(a)\|} \lesssim L_\psi \|\Delta a\| \leq L_\psi \|a\| \cdot \frac{\|\Delta a\|}{\|a\|}.$$

1.4 Интервальная арифметика

Исходной идеей создания интервальной арифметики является наблюдение о том, что всё в нашем мире неточно, и нам в реальности чаще всего приходится работать не с точными значениями величин, которые образуют основу классической «идеальной» математики, а с целыми диапазонами значений той или иной величины. Например, множество вещественных чисел, которые точно представляются в цифровых ЭВМ, конечно, и из-за присутствия округления каждое из этих чисел, в действительности, является представителем целого интервала значений обычной вещественной оси \mathbb{R} (см. Рис. 1.5–1.6).

Нельзя ли организовать операции и отношения между диапазонами-интервалами так, как это сделано для обычных точных значений? С

тем, чтобы можно было работать с ними, подобно обычным числам, опираясь на алгебраические преобразования, аналитические операции и т.п.? Ответ на эти вопросы положителен, хотя свойства получающейся «интервальной арифметики» оказываются во многом непохожими на привычные свойства операций с обычными числами.

Предположим, что нам даны переменные a и b , точные значения которых неизвестны, но мы знаем, что они могут находиться в интервалах $[\underline{a}, \bar{a}]$ и $[\underline{b}, \bar{b}]$ соответственно. Что можно сказать о значении суммы $a + b$?

Складывая почленно неравенства

$$\begin{aligned}\underline{a} &\leq a \leq \bar{a}, \\ \underline{b} &\leq b \leq \bar{b},\end{aligned}$$

получим

$$\underline{a} + \underline{b} \leq a + b \leq \bar{a} + \bar{b},$$

так что $a + b \in [\underline{a} + \underline{b}, \bar{a} + \bar{b}]$.

На аналогичный вопрос, связанный с областью значений разности $a - b$ можно ответить, складывая почленно неравенства

$$\begin{aligned}\underline{a} &\leq a \leq \bar{a}, \\ -\bar{b} &\leq -b \leq -\underline{b}.\end{aligned}$$

Имеем в результате $a - b \in [\underline{a} - \bar{b}, \bar{a} - \underline{b}]$.

Для умножения двух переменных $a \in [\underline{a}, \bar{a}]$ и $b \in [\underline{b}, \bar{b}]$ имеет место несколько более сложная оценка

$$a \cdot b \in [\min\{\underline{a}\underline{b}, \underline{a}\bar{b}, \bar{a}\underline{b}, \bar{a}\bar{b}\}, \max\{\underline{a}\underline{b}, \underline{a}\bar{b}, \bar{a}\underline{b}, \bar{a}\bar{b}\}].$$

Чтобы доказать её заметим, что функция $\phi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, задаваемая правилом $\phi(a, b) = a \cdot b$, будучи линейной по b при каждом фиксированном a , принимает минимальное и максимальное значения на концах интервала изменения переменной b . Это же верно и для экстремумов по $a \in [\underline{a}, \bar{a}]$ при любом фиксированном значении b . Наконец,

$$\begin{aligned}\min_{a \in [\underline{a}, \bar{a}], b \in [\underline{b}, \bar{b}]} \phi(a, b) &= \min_{a \in [\underline{a}, \bar{a}]} \min_{b \in [\underline{b}, \bar{b}]} \phi(a, b), \\ \max_{a \in [\underline{a}, \bar{a}], b \in [\underline{b}, \bar{b}]} \phi(a, b) &= \max_{a \in [\underline{a}, \bar{a}]} \max_{b \in [\underline{b}, \bar{b}]} \phi(a, b),\end{aligned}$$

т. е. взятие минимума по совокупности аргументов может быть заменено повторным минимумом, а взятие максимума по совокупности аргументов — повторным максимумом, причём в обоих случаях порядок экстремумов несуществен. Следовательно, для $a \in [\underline{a}, \bar{a}]$ и $b \in [\underline{b}, \bar{b}]$ в самом деле

$$\min\{\underline{a}\underline{b}, \underline{a}\bar{b}, \bar{a}\underline{b}, \bar{a}\bar{b}\} \leq a \cdot b \leq \max\{\underline{a}\underline{b}, \underline{a}\bar{b}, \bar{a}\underline{b}, \bar{a}\bar{b}\}, \quad (1.8)$$

и нетрудно видеть, что эта оценка достижима с обеих сторон.

Наконец, для частного двух ограниченных переменных несложно вывести оценки из неравенств для умножения и из того факта, что $a/b = a \cdot (1/b)$.

Проведённые выше рассуждения подсказывают идею — рассматривать интервалы вещественной оси как самостоятельные объекты, между которыми можно будет ввести свои собственные операции, отношения и т. п. Мы далее будем обозначать интервалы буквами жирного шрифта: \mathbf{a} , \mathbf{b} , \mathbf{c} , \dots , \mathbf{x} , \mathbf{y} , \mathbf{z} . Подчёркивание и надчёркивание — $\underline{\mathbf{a}}$ и $\bar{\mathbf{a}}$ — будут зарезервированы для обозначения нижнего и верхнего концов интервала, так что $\mathbf{a} = [\underline{\mathbf{a}}, \bar{\mathbf{a}}]$.

Рассмотрим множество всех вещественных интервалов $\mathbf{a} := [\underline{\mathbf{a}}, \bar{\mathbf{a}}] = \{a \in \mathbb{R} \mid \underline{\mathbf{a}} \leq a \leq \bar{\mathbf{a}}\}$, и бинарные операции — сложение, вычитание, умножение и деление — определим между ними «по представителям», т. е. в соответствии со следующим фундаментальным принципом:

$$\mathbf{a} \star \mathbf{b} := \{a \star b \mid a \in \mathbf{a}, b \in \mathbf{b}\} \quad (1.9)$$

для всех интервалов \mathbf{a} , \mathbf{b} , таких что выполнение точечной операции $a \star b$, $\star \in \{+, -, \cdot, /\}$, имеет смысл для любых $a \in \mathbf{a}$ и $b \in \mathbf{b}$. При этом вещественные числа отождествляются с интервалами нулевой ширины $[a, a]$, которые называются также *вырожденными интервалами*. Кроме того, через $(-\mathbf{a})$ условимся обозначать интервал $(-1) \cdot \mathbf{a}$.

Для интервальных арифметических операций развёрнутое определение, равносильное (1.9), как мы установили выше, задаётся следующими формулами:

$$\mathbf{a} + \mathbf{b} = [\underline{\mathbf{a}} + \underline{\mathbf{b}}, \bar{\mathbf{a}} + \bar{\mathbf{b}}], \quad (1.10)$$

$$\mathbf{a} - \mathbf{b} = [\underline{\mathbf{a}} - \bar{\mathbf{b}}, \bar{\mathbf{a}} - \underline{\mathbf{b}}], \quad (1.11)$$

$$\mathbf{a} \cdot \mathbf{b} = [\min\{\underline{\mathbf{a}}\underline{\mathbf{b}}, \underline{\mathbf{a}}\bar{\mathbf{b}}, \bar{\mathbf{a}}\underline{\mathbf{b}}, \bar{\mathbf{a}}\bar{\mathbf{b}}\}, \max\{\underline{\mathbf{a}}\underline{\mathbf{b}}, \underline{\mathbf{a}}\bar{\mathbf{b}}, \bar{\mathbf{a}}\underline{\mathbf{b}}, \bar{\mathbf{a}}\bar{\mathbf{b}}\}], \quad (1.12)$$

$$\mathbf{a}/\mathbf{b} = \mathbf{a} \cdot [1/\bar{\mathbf{b}}, 1/\underline{\mathbf{b}}] \quad \text{для } \mathbf{b} \not\equiv 0. \quad (1.13)$$

В частности, при умножении интервала на число полезно помнить следующее простое правило:

$$\mu \cdot \mathbf{a} = \begin{cases} [\mu \underline{\mathbf{a}}, \mu \bar{\mathbf{a}}], & \text{если } \mu \geq 0, \\ [\mu \bar{\mathbf{a}}, \mu \underline{\mathbf{a}}], & \text{если } \mu \leq 0. \end{cases} \quad (1.14)$$

Алгебраическая система $\langle \mathbb{IR}, +, -, \cdot, / \rangle$, образованная множеством всех вещественных интервалов $\mathbf{a} := [\underline{\mathbf{a}}, \bar{\mathbf{a}}] = \{x \in \mathbb{R} \mid \underline{\mathbf{a}} \leq x \leq \bar{\mathbf{a}}\}$ с бинарными операциями сложения, вычитания, умножения и деления, которые определены формулами (1.10)–(1.13), называется *классической интервальной арифметикой*. Эпитет «классическая» используется здесь потому, что существуют и другие интервальные арифметики, приспособленные для решения других задач.

Полезно выписать определение интервального умножения в виде так называемой таблицы Кэли, дающей представление результата операции в зависимости от различных комбинаций значений операндов. Для этого выделим в \mathbb{IR} следующие подмножества:

$$\begin{aligned} \mathcal{P} &:= \{\mathbf{a} \in \mathbb{IR} \mid \underline{\mathbf{a}} \geq 0 \text{ и } \bar{\mathbf{a}} \geq 0\} && \text{— неотрицательные интервалы,} \\ \mathcal{Z} &:= \{\mathbf{a} \in \mathbb{IR} \mid \underline{\mathbf{a}} \leq 0 \leq \bar{\mathbf{a}}\} && \text{— нульсодержащие интервалы,} \\ -\mathcal{P} &:= \{\mathbf{a} \in \mathbb{IR} \mid -\mathbf{a} \in \mathcal{P}\} && \text{— неположительные интервалы.} \end{aligned}$$

В целом $\mathbb{IR} = \mathcal{P} \cup \mathcal{Z} \cup (-\mathcal{P})$. Тогда интервальное умножение (1.12) может быть описано с помощью Табл. 1.1, особенно удобной при реализации этой операции на ЭВМ.

Именно по этой таблице реализовано интервальное умножение в подавляющем большинстве компьютерных систем, поддерживающих интервальную арифметику, так как в сравнении с исходными формулами такая реализация существенно более быстрая.

Алгебраические свойства классической интервальной арифметики существенно беднее, чем у поля вещественных чисел \mathbb{R} . В частности, особенностью интервальной арифметики является отсутствие дистрибутивности умножения относительно сложения: в общем случае

$$(\mathbf{a} + \mathbf{b})\mathbf{c} \neq \mathbf{ac} + \mathbf{bc}.$$

Например,

$$[1, 2] \cdot (1 - 1) = 0 \neq [-1, 1] = [1, 2] \cdot 1 - [1, 2] \cdot 1.$$

Таблица 1.1. Интервальное умножение

\cdot	$b \in \mathcal{P}$	$b \in \mathcal{Z}$	$b \in -\mathcal{P}$
$a \in \mathcal{P}$	$[\underline{a}b, \bar{a}\bar{b}]$	$[\bar{a}\underline{b}, \bar{a}\bar{b}]$	$[\bar{a}\underline{b}, \underline{a}\bar{b}]$
$a \in \mathcal{Z}$	$[\underline{a}\bar{b}, \bar{a}\bar{b}]$	$[\min\{\underline{a}\bar{b}, \bar{a}\underline{b}\}, \max\{\underline{a}\underline{b}, \bar{a}\bar{b}\}]$	$[\bar{a}\underline{b}, \underline{a}\bar{b}]$
$a \in -\mathcal{P}$	$[\underline{a}\bar{b}, \bar{a}\underline{b}]$	$[\underline{a}\bar{b}, \underline{a}\underline{b}]$	$[\bar{a}\bar{b}, \underline{a}\bar{b}]$

Тем не менее, имеет место более слабое свойство

$$a(b + c) \subseteq ab + ac \tag{1.15}$$

называемое *субдистрибутивностью* умножения относительно сложения. В ряде частных случаев дистрибутивность всё-таки выполняется:

$$a(b + c) = ab + ac, \quad \text{если } a \text{ — вещественное число,} \tag{1.16}$$

$$a(b + c) = ab + ac, \quad \text{если } b, c \geq 0 \text{ или } b, c \leq 0. \tag{1.17}$$

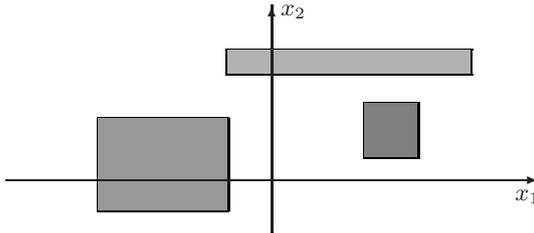


Рис. 1.3. Интервальные векторы-брусы в \mathbb{R}^2 .

Интервальный вектор — это упорядоченный кортеж из интервалов, расположенный вертикально (вектор-столбец) или горизонтально

(вектор-строка). Таким образом, если $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$ — некоторые интервалы, то

$$\mathbf{a} = \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_n \end{pmatrix} \text{ — это интервальный вектор-столбец,}$$

а

$$\mathbf{a} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n) \text{ — это интервальная вектор-строка.}$$

Множество интервальных векторов, компоненты которых принадлежат \mathbb{IR} , мы будем обозначать через \mathbb{IR}^n . При этом нулевые векторы, т. е. такие, все компоненты которых суть нули, мы традиционно обозначаем через «0».

Введём также важное понятие интервальной оболочки множества. Если S — непустое ограниченное множество в \mathbb{R}^n или $\mathbb{R}^{m \times n}$, то его *интервальной оболочкой* $\square S$ называется наименьший по включению интервальный вектор (или матрица), содержащий S . Нетрудно понять, что это определение равносильно такому: интервальная оболочка множества S — это пересечение всех интервальных векторов, содержащих S , т. е.

$$\square S = \cap \{ \mathbf{a} \in \mathbb{IR}^n \mid \mathbf{a} \supseteq S \}.$$

Интервальная оболочка — это интервальный объект, наилучшим образом приближающий извне (т. е. объемлющий) рассматриваемое множество, и компоненты $\square S$ являются проекциями множества S на координатные оси пространства.

Сумма (разность) двух интервальных матриц одинакового размера есть интервальная матрица того же размера, образованная поэлементными суммами (разностями) операндов.

Если $\mathbf{A} = (\mathbf{a}_{ij}) \in \mathbb{IR}^{m \times l}$ и $\mathbf{B} = (\mathbf{b}_{ij}) \in \mathbb{IR}^{l \times n}$, то произведение матриц \mathbf{A} и \mathbf{B} есть матрица $\mathbf{C} = (\mathbf{c}_{ij}) \in \mathbb{IR}^{m \times n}$, такая что

$$\mathbf{c}_{ij} := \sum_{k=1}^l \mathbf{a}_{ik} \mathbf{b}_{kj}.$$

Нетрудно показать, что для операций между матрицами выполняется соотношение

$$\mathbf{A} \star \mathbf{B} = \square \{ \mathbf{A} \star \mathbf{B} \mid \mathbf{A} \in \mathbf{A}, \mathbf{B} \in \mathbf{B} \}, \quad \star \in \{ +, -, \cdot \}, \quad (1.18)$$

где \square — интервальная оболочка множества, наименьший по включению интервальный вектор-брус, который содержит его.

1.5 Интервальные расширения функций

Пусть $f : \mathbb{R} \rightarrow \mathbb{R}$ — некоторая функция. Если мы рассматриваем интервалы в виде самостоятельных объектов, то что следует понимать под значением функции от интервала? Естественно считать, что

$$f(\mathbf{x}) = \{f(x) \mid x \in \mathbf{x}\}.$$

Задача об определении области значений функции на том или ином подмножестве области её определения, эквивалентная задаче оптимизации, в интервальном анализе принимает специфическую форму задачи о вычислении так называемого *интервального расширения функции*.

Определение 1.5.1 Пусть D — непустое подмножество пространства \mathbb{R}^n . Интервальная функция $\mathbf{f} : \mathbb{I}D \rightarrow \mathbb{I}\mathbb{R}^m$ называется интервальным продолжением точечной функции $f : D \rightarrow \mathbb{R}^m$, если $\mathbf{f}(x) = f(x)$ для всех $x \in D$.

Определение 1.5.2 Пусть D — непустое подмножество пространства \mathbb{R}^n . Интервальная функция $\mathbf{f} : \mathbb{I}D \rightarrow \mathbb{I}\mathbb{R}^m$ называется интервальным расширением точечной функции $f : D \rightarrow \mathbb{R}^m$, если

- 1) $\mathbf{f}(x)$ — интервальное продолжение $f(x)$,
- 2) $\mathbf{f}(x)$ монотонна по включению, т. е.

$$\mathbf{x}' \subseteq \mathbf{x}'' \Rightarrow \mathbf{f}(\mathbf{x}') \subseteq \mathbf{f}(\mathbf{x}'') \text{ на } \mathbb{I}D.$$

Таким образом, если $\mathbf{f}(x)$ — интервальное расширение функции $f(x)$, то для области значений f на бресе $\mathbf{X} \subset D$ мы получаем следующую внешнюю (с помощью объемлющего множества) оценку:

$$\{f(x) \mid x \in \mathbf{X}\} = \bigcup_{x \in \mathbf{X}} f(x) = \bigcup_{x \in \mathbf{X}} \mathbf{f}(x) \subseteq \mathbf{f}(\mathbf{X}).$$

Эффективное построение интервальных расширений функций — это важнейшая задача интервального анализа, поиски различных решений которой продолжаются и в настоящее время. Уместно привести в

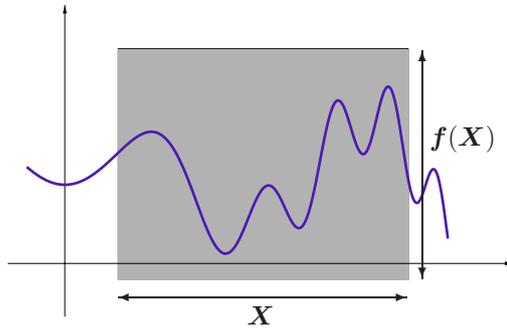


Рис. 1.4. Интервальное расширение функции даёт внешнюю оценку её области значений

рамках нашего беглого обзора некоторые общезначимые результаты в этом направлении. Первый из них часто называют «основной теоремой интервальной арифметики»:

Теорема 1.5.1 Если для рациональной функции $f(x) = f(x_1, x_2, \dots, x_n)$ на бруске $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{IR}^n$ определён результат $\mathbf{f}_i(\mathbf{x})$ подстановки вместо её аргументов интервалов их изменения x_1, x_2, \dots, x_n и выполнения всех действий над ними по правилам интервальной арифметики, то

$$\{f(x) \mid x \in \mathbf{x}\} \subseteq \mathbf{f}(\mathbf{x}),$$

т. е. $\mathbf{f}(\mathbf{x})$ содержит множество значений функции $f(x)$ на \mathbf{x} .

Нетрудно понять, что по отношению к рациональной функции $f(x)$ интервальная функция $\mathbf{f}_i(\mathbf{x})$, о которой идёт речь в Теореме 1.5.1, является интервальным расширением. Оно называется *естественным интервальным расширением* и вычисляется совершенно элементарно.

Пример 1.5.1 Для функции $f(x) = x/(x+1)$ на интервале $[1, 3]$ область значений в соответствии с результатом Теоремы 1.5.1 можно оценить извне как

$$\frac{[1, 3]}{[1, 3] + 1} = \frac{[1, 3]}{[2, 4]} = \left[\frac{1}{4}, \frac{3}{2}\right]. \quad (1.19)$$

Но если предварительно переписать выражение для функции в виде

$$f(x) = \frac{1}{1 + 1/x},$$

разделив числитель и знаменатель дроби на $x \neq 0$, то интервальное оценивание даст уже результат

$$\frac{1}{1 + 1/[1, 3]} = \frac{1}{[\frac{4}{3}, 2]} = [\frac{1}{2}, \frac{3}{4}].$$

Он более узок (т.е. более точен), чем (1.19), и совпадает к тому же с областью значений. Как видим качество интервального оценивания существенно зависит от вида выражения. ■

Использование естественного интервального расширения подчас даёт весьма грубые оценки областей значений функций, в связи с чем получили развитие более совершенные способы (формы) нахождения интервальных расширений. Одна из наиболее популярных — так называемая *центрированная форма*:

$$\mathbf{f}_c(\mathbf{x}, \tilde{x}) = f(\tilde{x}) + \sum_{i=1}^n \mathbf{g}_i(\mathbf{x}, \tilde{x})(\mathbf{x}_i - \tilde{x}_i),$$

где $\tilde{x} = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$ — некоторая фиксированная точка, называемая «центром»,

$\mathbf{g}_i(\mathbf{x}, \tilde{x})$ — интервальные расширения некоторых функций $g_i(x, \tilde{x})$, которые строятся по f и зависят в общем случае как от \tilde{x} , так и от \mathbf{x} .

В выписанном выше выражении $\mathbf{g}_i(\mathbf{x}, \tilde{x})$ могут быть внешними оценками коэффициентов наклона функции f на рассматриваемой области определения, взятыми относительно точки \tilde{x} , или же внешними интервальными оценками областей значений производных $\partial f(x)/\partial x_i$ на \mathbf{x} . В последнем случае точка \tilde{x} никак не используется, а интервальная функция \mathbf{f}_c называется дифференциальной центрированной формой интервального расширения.²

Пример 1.5.2 Для оценивания функции $f(x) = x/(x+1)$ на интервале $x = [1, 3]$ применим дифференциальную центрированную форму.

Так как

$$f'(x) = \frac{1}{(x+1)^2},$$

²По отношению к ней часто используют также термин «среднезначная форма», поскольку она может быть выведена из известной теоремы Лагранжа о среднем.

то интервальная оценка производной на заданном интервале области определения есть

$$\frac{1}{([1, 3] + 1)^2} = \left[\frac{1}{16}, \frac{1}{4} \right]$$

Поэтому если в качестве центра разложения взять середину интервала $\text{mid } \mathbf{x} = 2$, то

$$\begin{aligned} f(\text{mid } \mathbf{x}) + f'(\mathbf{x})(\mathbf{x} - \text{mid } \mathbf{x}) &= \frac{2}{3} + \left[\frac{1}{16}, \frac{1}{4} \right] \cdot [-1, 1] \\ &= \frac{2}{3} + \left[-\frac{1}{4}, \frac{1}{4} \right] = \left[\frac{5}{12}, \frac{11}{12} \right]. \end{aligned}$$

Как видим, этот результат значительно точнее естественного интервального расширения (1.19). ■

За дальнейшей информацией мы отсылаем заинтересованного читателя к книгам [1, 24, 28, 29], развёрнуто излагающим построение интервальных расширений функций. Важно отметить, что точность интервального оценивания при использовании любой из форм интервального расширения критическим образом зависит от ширины бруса оценивания. Если обозначить через $f(\mathbf{x})$ точную область значений целевой функции на \mathbf{x} , т. е. $f(\mathbf{x}) = \{f(x) \mid x \in \mathbf{x}\}$, то для естественного интервального расширения липшицевых функций имеет место неравенство

$$\text{dist}(\mathbf{f}_i(\mathbf{x}), f(\mathbf{x})) \leq C \|\text{wid } \mathbf{x}\| \quad (1.20)$$

с некоторой константой C , и этот факт обычно выражают словами «естественное интервальное расширение имеет первый порядок точности». Для центрированной формы верно соотношение

$$\text{dist}(\mathbf{f}_c(\mathbf{x}, \tilde{x}), f(\mathbf{x})) \leq 2(\text{wid } \mathbf{g}(\mathbf{x}, \tilde{x}))^\top |\mathbf{x} - \tilde{x}|, \quad (1.21)$$

где $\mathbf{g}(\mathbf{x}, \tilde{x}) = (\mathbf{g}_1(\mathbf{x}, \tilde{x}), \mathbf{g}_2(\mathbf{x}, \tilde{x}), \dots, \mathbf{g}_n(\mathbf{x}, \tilde{x}))$. В случае, когда интервальные оценки для функций $\mathbf{g}_i(\mathbf{x}, \tilde{x})$ находятся с первым порядком точности, общий порядок точности центрированной формы согласно (1.21) будет уже вторым. Вывод этих оценок заинтересованный читатель может найти, к примеру, в [24, 29].

Интервальные оценки областей значений функций, которые находятся с помощью интервальных расширений, оказываются полезными в самых различных вопросах вычислительной математики. В частности, с помощью интервального языка очень элегантно записываются остаточные члены различных приближённых формул. В качестве

двух содержательных примеров применения интервальных расширений функций мы рассмотрим решение уравнений и оценку константы Липшица для функций.

1.6 Элементы конструктивной математики

«Конструктивная математика» — это неформальное название той части современной математики, тех математических дисциплин, — теории алгоритмов, теории сложности вычислений, и ряда других — в которых главным объектом изучения являются процессы построения тех или иных математических объектов. Оформление конструктивной математики в отдельную ветвь общего математического дерева произошло на рубеже XIX и XX веков под влиянием обнаруженных к тому времени парадоксов теории множеств. Эти парадоксы заставили критически переосмыслить существовавшие в математике способы рассуждений и само понятие «существования» для математических объектов. Создание основ конструктивного направления в математике связано, прежде всего, с деятельностью Л.Э.Я. Брауэра и развиваемым им «интуиционизмом».

В частности, теория алгоритмов и рекурсивных функций — это математическая дисциплина, исследующая конструктивные свойства различных математических объектов. Её основные понятия — это *алгоритм*, *конструктивный объект*, *вычислимость*, *разрешимость* и др.

Алгоритм — это конечная последовательность инструкций, записанных на некотором языке и определяющих процесс переработки исходных данных в искомые результаты (ответ решаемой задачи и т.п.). Алгоритм принципиально конечен и определяет собой конечный процесс. Далее, *конструктивным объектом* называется объект, который может быть построен с помощью некоторой конечной последовательности действий над каким-то конечным алфавитом. Таковы, например, рациональные числа. Строго говоря, конструктивные объекты и только они могут быть получены в качестве ответов при решении задачи на реальных цифровых ЭВМ с конечными быстродействием и объёмом памяти.

В частности, конечными машинами являются широко распространенные ныне электронные цифровые вычислительные машины: они способны представлять, по сути дела, только конечные множества чисел. Таким образом, обречены на неудачу любые попытки использо-

вать их для выполнения арифметических абсолютно точных операций над числовыми полями \mathbb{R} и \mathbb{C} , которые являются бесконечными (и даже непрерывными) множествами, большинство элементов которых не представимы в цифровых ЭВМ.

Оказывается, что значительная часть объектов, с которыми работают современная математика и её приложения, не являются конструктивными. В частности, неконструктивным является традиционное понятие вещественного числа, подразумевающее бесконечную процедуру определения всех знаков его десятичного разложения (которое в общем случае неперiodично). Факт неконструктивности вещественных чисел может быть обоснован строго математически (см. [32]), и он указывает на принципиальные границы возможностей алгоритмического подхода и ЭВМ в деле решения задач математического анализа.

Тем не менее, и в этом океане неконструктивности имеет смысл выделить объекты, которые могут быть «достаточно хорошо» приближены конструктивными объектами. На этом пути мы приходим к понятию *вычислимого вещественного числа* [32, 22]³: вещественное число α называется вычислимым, если существует алгоритм, дающий по всякому натуральному числу n рациональное приближение к α с погрешностью $\frac{1}{n}$. Множество всех вычисляемых вещественных чисел образует *вычисляемый континуум*. Соответственно, *вычисляемая вещественная функция* определяется как отображение из вычислимого континуума в вычисляемый континуум, задаваемая алгоритмом преобразования программы аргумента в программу значений.

Важно помнить, что и вычисляемое вещественное число, и вычисляемая функция — это уже не конструктивные объекты. Но, как выясняется, даже ценой ослабления наших требований к конструктивности нельзя вполне преодолеть принципиальные алгоритмические трудности, связанные с задачей решения уравнений. Для вычисляемых вещественных чисел и функций ряд традиционных постановок задач оказывается *алгоритмически неразрешимыми* в том смысле, что построение общих алгоритмов их решения принципиально невозможно.

Например, алгоритмически неразрешимыми являются задачи

- 1) распознавания равно нулю или нет произвольное вычисляемое вещественное число [31, 32, 33], распознавания равенства двух вычисляемых вещественных чисел [22, 25, 31, 32];

³Совершенно аналогичным является определение *конструктивного вещественного числа* у Б.А. Кушнера [31].

- 2) нахождения для каждой совместной системы линейных уравнений над полем конструктивных вещественных чисел какого-либо ее решения [31, 33];
- 3) нахождения нулей всякой непрерывной кусочно-линейной знакопеременной функции [33].

Приведённые выше результаты задают, как нам представляется, ту абсолютную и совершенно объективную мерку (в отличие от субъективных пристрастий), с которой мы должны подходить к оценке трудоёмкости тех или иных вычислительных методов. Получается, что необходимость переформулировки задачи решения уравнений и систем уравнений связана ещё и с тем, что в традиционной постановке эти задачи оказываются алгоритмически неразрешимыми! На фоне этого мрачного факта наличие даже экспоненциально трудного алгоритма с небольшим основанием «одноэтажной» экспоненты в оценке сложности (вроде 2^n) можно рассматривать как вполне приемлемый вариант разрешимости задачи. Именно это имеет место в ситуации с вычислением вращения векторного поля (степени отображения).

Вычислительная математика тесно примыкает к конструктивной, хотя их цели и методы существенно разнятся.

1.7 Сложность задач и трудоёмкость алгоритмов

Как правило, нас удовлетворит не всякий процесс решения поставленной задачи, а лишь только тот, который выполним за практически приемлемое время. Соответственно, помимо алгоритмической разрешимости задач огромную роль играет трудоёмкость тех или иных алгоритмов для их решения.

Например, множество вещественных чисел, точно представимых в цифровых ЭВМ в формате «с плавающей точкой» согласно стандарту IEEE 754/854, является конечным, и потому мы можем найти, скажем, приближённые значения корней полинома (или убедиться в их отсутствии) за конечное время, просто перебрав все эти машинные числа и вычисляя в них значения полинома. Но, будучи принципиально выполнимым, такой алгоритм требует непомерных вычислительных затрат и для практики бесполезен.

Естественно измерять трудоёмкость алгоритма количеством «элементарных операций», требуемых для его исполнения. Следует лишь иметь в виду, что эти операции могут быть весьма различными. Скажем, сложение и умножение двух чисел «с плавающей точкой» требуют для своего выполнения разного количества тактов современных процессоров и, соответственно, разного времени. До определённой степени эти различия можно игнорировать и оперировать понятием усреднённой арифметической операции.

Большую роль играет также объём данных, подаваемых на вход алгоритма. К примеру, входными данными могут быть небольшие целые числа, а могут и рациональные дроби с внушительными числителями и знаменателями. Ясно, что переработка больших объёмов данных должна потребовать больших трудозатрат от алгоритма, так что имеет смысл сложность исполнения алгоритма в каждом конкретном случае отнести к сложности представления входных данных алгоритма.⁴

На качественном уровне полезно различать *полиномиальную трудоёмкость* и *экспоненциальную трудоёмкость*. Говорят, что некоторый алгоритм имеет полиномиальную трудоёмкость, если сложность его выполнения не превосходит значений некоторого полинома от длины входных данных. Напротив, про некоторый алгоритм говорят, что он имеет экспоненциальную трудоёмкость, если сложность его выполнения превосходит значения любого полинома от длины подаваемых ему на вход данных.

Получение оценок трудоёмкости задач является непростым делом. Если какой-то алгоритм решает поставленную задачу, то, очевидно, его трудоёмкость может служить верхней оценкой сложности решения этой задачи. Но вот получение нижних оценок сложности решения задач является чрезвычайно трудным. В явном виде такие нижние оценки найдены лишь для очень небольшого круга задач, которые имеют, скорее, теоретическое значение. В этих условиях широкое распространение получила альтернативная теория сложности, в основе которой лежат понятие сводимости задач друг к другу и вытекающее из него понятие эквивалентности задач по трудоёмкости.

Наибольшее распространение получило *полиномиальное сведение* одних задач к другим, под которым понимается такое преобразование одной задачи к другой, что данные и ответ исходной задачи перево-

⁴В связи с этим получили распространение также относительные единицы измерения трудоёмкости алгоритмов — через количество вычислений функции, правой части уравнения и т. п.

дятся в данные и ответ для другой, а трудоёмкость этого преобразования не превышает значений некоторого полинома от размера исходной задачи. Взаимная полиномиальная сводимость двух задач друг другу является отношением эквивалентности.

Этим рассуждениям можно придать более строгую форму, что приводит к так называемой теории NP-трудности, получившей существенное развитие в последние десятилетия. Её основным понятием является понятие NP-трудной задачи (универсальной переборной задачи) [10]

Таким образом, теория NP-полноты не отвечает напрямую на вопрос о трудоёмкости решения тех или иных задач, но позволяет утверждать, что некоторые задачи «столь же трудны», как и другие известные задачи. Нередко знание уже этого одного факта бывает существенным для ориентировки создателям вычислительных технологий решения конкретных задач. Если известно, к примеру, что некоторая задача не проще, чем известные «переборные» задачи, которые, по видимому, не могут быть решены лучше, чем полным перебором всех возможных вариантов, то имеет смысл и для рассматриваемой задачи не стесняться конструирования алгоритмов «переборного» типа, имеющих экспоненциальную трудоёмкость.

Именно такова ситуация с некоторыми задачами, которые возникают в вычислительной математике. К примеру, оценивание разброса решений систем линейных или нелинейных уравнений при варьировании параметров этих систем в самом общем случае, когда мы не ограничиваем себя величиной возмущений, является NP-трудной задачей. В частности, таковы интервальные системы уравнений.

1.8 Доказательные вычисления на ЭВМ

Термин «доказательные вычисления» был введён в 70-е годы XX века советским математиком К.И. Бабенко для обозначения вычислений, результат которых имеет такой же статус достоверности, как и результаты «чистой математики», полученные с помощью традиционных доказательств. В книге [3], где доказательным вычислениям посвящён отдельный параграф, можно прочитать: «Под *доказательными вычислениями* в анализе мы понимаем такие целенаправленные вычисления на ЭВМ, комбинируемые с аналитическими исследованиями, которые приводят к строгому установлению новых фактов (теорем)». В отношении задач, где ответом являются числа (набор чисел, вектор или

матрица и т.п.) доказательность означает свойство гарантированности этих числовых ответов.⁵ К примеру, если мы находим число π , то доказательным ответом может быть установление гарантированного неравенства $\pi > 3.1415926$ или $\pi \geq 3.1415926$.

Основная трудность, с которой сталкиваются при проведении доказательных вычислений на современных цифровых ЭВМ, вытекает из невозможности адекватно отобразить непрерывную числовую ось \mathbb{R} в виде множества машинно представимых чисел. Таковых может быть лишь конечное число (либо потенциально счётное), тогда как вещественная ось \mathbb{R} является непрерывным континуумом. Как следствие, типичное вещественное число не представимо точно в цифровой ЭВМ с конечной разрядной сеткой.



Рис. 1.5. Типичное вещественное число не представимо точно в цифровой ЭВМ с конечной разрядной сеткой

Ситуация, в действительности, ещё более серьёзна, так как неизбежными ошибками, как правило, сопровождаются ввод данных в ЭВМ и выполнение с ними арифметических операций. Хотя эти ошибки могут быть очень малы, но, накапливаясь, они способны существенно исказить ответ к решаемой задаче. Встаёт нетривиальная проблема их учёта в процессе счёта на ЭВМ ...



Рис. 1.6. Интервальное решение проблемы представления вещественных чисел в цифровой ЭВМ

Одним из средств доказательных вычислений на ЭВМ служит интервальная арифметика и, более общо, методы интервального анализа. В частности, вещественное число x в общем случае наиболее корректно

⁵Термин «доказательные вычисления на ЭВМ» является хорошим русским эквивалентом таких распространённых английских оборотов как *verified computation*, *verification numerics* и др.

представляется в цифровых ЭВМ интервалом, левый конец которого — наибольшее машинно-представимое число, не превосходящее x , а правый — наименьшее машинно-представимое число, не меньшее x . Далее с получающимися интервалами можно выполнять операции по правилам интервальной арифметики, рассмотренным в §1.4.

Концы интервалов, получающихся при расчётах по формулам (1.10)–(1.13), также могут оказаться вещественными числами, не представимыми в ЭВМ. В этом случае для обеспечения доказательности вычислений имеет смысл несколько расширить полученный интервал до ближайшего объемлющего его интервала с машинно-представимыми концами. Подобная версия интервальной арифметики называется *машинной интервальной арифметикой* с направленным округлением (см. Рис. 1.7).

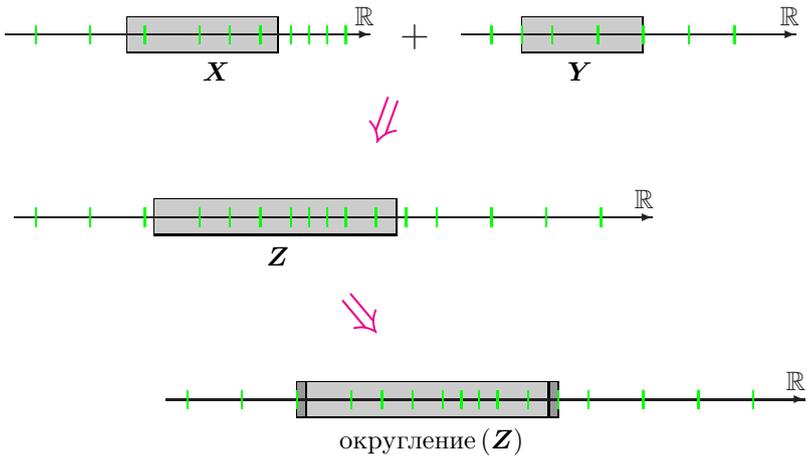


Рис. 1.7. Машинная интервальная арифметика с внешним направленным округлением

Существует несколько подходов к организации доказательных вычислений на ЭВМ, из которых наиболее известными являются *пошаговый способ* оценки ошибок и *апостериорное оценивание*.

В пошаговом способе доказательных вычислений мы разбиваем алгоритм вычисления решения на «элементарные шаги», оцениваем погрешности на каждом шаге вычислений. «Элементарными шагами» здесь могут быть как отдельные арифметические и логические опера-

ции, так и целые их последовательности, слагающиеся в крупные блоки алгоритма. При этом полная погрешность получается из погрешностей отдельных «элементарных шагов» по правилам исчисления из §1.2. Очевидный недостаток такого способа организации оценки погрешностей состоит в том, что мы неявно привязываемся к конкретному алгоритму вычисления решения. При этом качество оценок, получаемых с помощью пошаговой парадигмы, существенно зависит от алгоритма, и «хороший» в обычном смысле алгоритм не обязательно хорош при оценивании погрешностей.

При оценивании погрешностей простых «элементарных шагов» алгоритмов с помощью таких несложных средств как классическая интервальная арифметика, получаемые оценки, как правило, отличаются невысоким качеством. Но изощёренные варианты пошагового способа оценки погрешностей могут показывать вполне удовлетворительные результаты даже для довольно сложных задач. Таковы, к примеру, вычислительные алгоритмы для решения систем линейных алгебраических уравнений, развиваемые в [30].

Напротив, при апостериорном оценивании погрешности мы оцениваем погрешность окончательного результата уже *после* его получения. Иными словами, мы разделяем способ получения двусторонней оценки решения и установление её доказательности. Ниже в Главе 4 мы приведём примеры конкретных алгоритмов апостериорного оценивания для доказательного решения некоторых популярных математических задач.

Литература к главе 1

Основная

- [1] АЛЕФЕЛЬД Г., ХЕРЦБЕРГЕР Ю. *Введение в интервальные вычисления*. – Москва: Мир, 1987.
- [2] БАРАХНИН В.Б., ШАПЕЕВ В.П. *Введение в численный анализ*. – Санкт-Петербург–Москва–Краснодар: Лань, 2005.
- [3] БАБЕНКО К.И. *Основы численного анализа*. – Москва: Наука, 1986.
- [4] БАУЭР Ф.Л., ГООЗ Г. *Информатика. В 2-х ч.* – Москва: Мир, 1990.
- [5] БАХВАЛОВ Н.С., ЖИДКОВ Н.П., КОБЕЛЬКОВ Г.М. *Численные методы*. – Москва: Бином, 2003, а также другие издания этой книги.
- [6] БАХВАЛОВ Н.С., КОРНЕВ А.А., ЧИЖОНКОВ Е.В. *Численные методы. Решения задач и упражнения*. – Москва: Дрофа, 2008.

- [7] БЕРЕЗИН И.С., ЖИДКОВ Н.П. *Методы вычислений. Т. 1–2.* – Москва: Наука, 1966.
- [8] ВЕРЖВИЦКИЙ В.М. *Численные методы. Части 1–2.* – Москва: «Оникс 21 век», 2005.
- [9] ВОЛКОВ Е.А. *Численные методы.* – Москва: Наука, 1987.
- [10] ГЭРИ М., ДЖОНСОН Д. *Вычислительные машины и труднорешаемые задачи.* – Москва: Мир, 1982.
- [11] ДЕМИДОВИЧ Б.П., МАРОН А.А. *Основы вычислительной математики.* – Москва: Наука, 1970.
- [12] КАЛИТКИН Н.Н. *Численные методы.* – Москва: Наука, 1978.
- [13] КРЫЛОВ В.И., БОБКОВ В.В., МОНАСТЫРНЫЙ П.И. *Вычислительные методы. Т. 1–2.* – Москва: Наука, 1976.
- [14] КУНЦ К.С. *Численный анализ.* – Киев: Техника, 1964.
- [15] КУНЦМАН Ж. *Численные методы.* – Москва: Наука, 1979.
- [16] МАЦОКИН А.М., СОРОКИН С.Б. *Численные методы. Часть 1. Численный анализ.* – Новосибирск: НГУ, 2006.
- [17] МЕНЬШИКОВ Г.Г. *Локализуемые вычисления. Конспект лекций.* – Санкт-Петербург: СПбГУ, Факультет прикладной математики–процессов управления, 2003.
- [18] МИНЬКОВ С.Л., МИНЬКОВ Л.Л. *Основы численных методов.* – Томск: Издательство научно-технической литературы, 2005.
- [19] РАЙС ДЖ. *Матричные вычисления и математическое обеспечение.* – Москва: Мир, 1984.
- [20] САМАРСКИЙ А.А., ГУЛИН А.В. *Численные методы.* – Москва: Наука, 1989.
- [21] ТЫРТЫШНИКОВ Е.Е. *Методы численного анализа.* – Москва: Академия, 2007.
- [22] УСПЕНСКИЙ В.А., СЕМЁНОВ А.Л. *Теория алгоритмов: основные открытия и приложения.* – Москва: Наука, 1987.
- [23] ХАНСЕН Э., УОЛСТЕР ДЖ.У. *Глобальная оптимизация с помощью методов интервального анализа.* – Москва-Ижевск: Издательство «РХД», 2012.
- [24] ШАРЫЙ С.П. *Конечномерный интервальный анализ.* – Электронная книга, 2010 (см. <http://www.nsc.ru/interval/Library/InteBooks>)
- [25] ABERNETHY O. *Precise numerical methods using C++.* – San Diego: Academic Press, 1998.
- [26] GOLDBERG D. What every computer scientist should know about floating point arithmetic // *ACM Computing Surveys.* – 1991. – Vol. 23, No. 1. – P. 5–48.
- [27] KREINOVICH V., LAKEYEV A.V., ROHN J., KAHL P. *Computational complexity and feasibility of data processing and interval computations.* – Dordrecht: Kluwer, 1997.
- [28] MOORE R.E., KEARFOTT R.B., CLOUD M. *Introduction to interval analysis.* – Philadelphia: SIAM, 2009.

- [29] NEUMAIER A. *Interval methods for systems of equations*. – Cambridge: Cambridge University Press, 1990.

Дополнительная

- [30] Годунов С.К., Антонов А.Г., Кирилюк О.Г., Костин В.И. *Гарантированная точность решения систем линейных уравнений в евклидовых пространствах*. – Новосибирск: Наука, 1988 и 1992.
- [31] КУШНЕР Б.А. *Лекции по конструктивному математическому анализу*. – Москва: Наука, 1973.
- [32] МАРТИН-ЛЁФ П. *Очерки по конструктивной математике*. – Москва: Наука, 1975.
- [33] *Математический Энциклопедический Словарь*. – Москва: Большая Российская Энциклопедия, 1995.
- [34] *IEEE Std 754-1985. IEEE Standard for Binary Floating-Point Arithmetic*. – New York: IEEE, 1985.

Глава 2

Численные методы анализа

2.1 Введение

Под численными методами анализа обычно понимаются вычислительные методы решения ряда задач, имеющих происхождение в классическом математическом анализе. Традиционно сюда относят задачи интерполирования и приближения функций, задачи численного нахождения производных и интегралов, а также задачу суммирования рядов. Кроме того, численные методы анализа охватывают задачу вычисления значений функций, которая относительно проста для функций, явно задаваемых несложными арифметическими выражениями, но становится нетривиальной в случае, когда функция задаётся неявно или с помощью операций, выводящих за пределы конечного набора элементарных арифметических действий.

В нашем курсе мы будем заниматься первыми четырьмя из перечисленных выше задач, и рассмотрим сначала задачи интерполирования и приближения функций.

Задачи интерполирования¹ и приближения функций являются тесно связанными друг с другом задачами, которые укладываются в рамки следующей единой неформальной схемы. Пусть дана функция $f(x)$, принадлежащая некоторому классу функций \mathcal{F} , и пусть также задан

¹Наряду с термином «интерполирование» в равной мере используется его синоним «интерполяция».

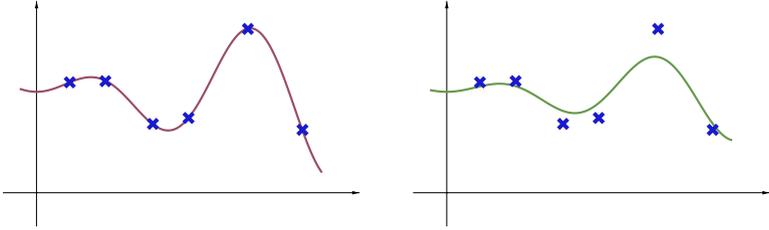


Рис. 2.1. Различие задач интерполяции и приближения функций.

класс функций \mathcal{G} . Требуется найти функцию $g(x)$ из \mathcal{G} , которая в определённом заранее смысле «достаточно близка» (или даже «наиболее близка») к данной функции $f(x)$. В зависимости от смысла, который вкладывается в понятие «близости» функций, в зависимости от того, какие именно функции образуют классы \mathcal{F} и \mathcal{G} , здесь могут получаться различные конкретные постановки задач. При этом полезно наделять рассматриваемые классы функций дополнительной структурой, например, считать, что они являются линейными векторными пространствами с нормой и т. п. Наконец, часто имеет место включение $\mathcal{G} \subset \mathcal{F}$.

При уточнении понятий «близости» функций и «отклонения» одной функции от другой обычно вводят множество значений аргументов X , на котором значения этих функций сравниваются друг с другом. X — это подмножество области определения функций, которое может совпадать со всей этой областью определения, но может также быть его небольшой частью, скажем, конечным набором точек. В последнем случае говорят о *дискретной* задаче приближения.

Задача интерполирования получается из приведённой выше общей формулировки в случае, когда «близость» означает совпадение функций f и g на некотором дискретном множестве точек x_0, x_1, \dots, x_n из пересечения их областей определения. От функции f при этом требуются лишь значения на этом множестве точек, и потому при постановке задачи интерполяции она сама часто даже не фигурирует. Вместо f обычно задаются лишь её значения y_0, y_1, \dots, y_n в точках x_0, x_1, \dots, x_n соответственно.

Задача приближения функций (называемая также задачей *аппроксимации функций*) является частным случаем общей формулировки, в котором «близость» понимается как малое отклонение значений функ-

ций f от g на подмножестве X из их области определения. Если X совпадает со всей областью определения, то удобно рассматривать эту «близость» или «отклонение» в терминах какого-то абстрактного расстояния (метрики), заданного на пересечении классов функций \mathcal{F} и \mathcal{G} (см. подробности в §2.10а). При этом в отличие от задачи интерполяции точное равенство функции g заданным значениям не требуется; ситуация иллюстрируется Рис. 2.1.

Напомним, что на множестве Y , образованном элементами произвольной природы, *расстоянием* (называемым также *метрикой*) называется определённая на декартовом произведении $Y \times Y$ функция dist с неотрицательными вещественными значениями, удовлетворяющая для любых $f, g, h \in Y$ следующим условиям:

- (1) $\text{dist}(f, g) = 0$ тогда и только тогда, когда $f = g$,
- (2) $\text{dist}(f, g) = \text{dist}(g, f)$ — симметричность,
- (3) $\text{dist}(f, h) \leq \text{dist}(f, g) + \text{dist}(g, h)$ — неравенство треугольника.

Разнообразные способы определения расстояния между функциями, возникающие в практике математического моделирования, приводят к различным математическим задачам приближения. Например, популярны равномерное (чебышёвское) отклонение функций друг от друга, которое определяется как

$$\max_{x \in [a, b]} |f(x) - g(x)|, \quad (2.1)$$

или интегральное отклонение функций на $[a, b]$, определяемое как

$$\int_a^b |f(x) - g(x)| dx. \quad (2.2)$$

В §2.10г мы рассмотрим также задачу среднеквадратичного приближения функций, в которой отклонение функций $f(x)$ и $g(x)$ друг от друга на интервале $[a, b]$ полагается равным

$$\left(\int_a^b (f(x) - g(x))^2 dx \right)^{1/2}. \quad (2.3)$$

Кроме перечисленных выше применяются также другие расстояния между функциями. Отметим, что расстояния (2.1)–(2.3) не вполне эквивалентны друг другу в том смысле, что сходимость последовательности функций к какому-то пределу относительно одного из этих расстояний не обязательно влечёт сходимость относительно другого.

Отметим, что в задачах дискретного приближения функций «отклонение» или «близость» хорошо охватываются понятием *псевдорасстояния* (псевдометрики), которое определяется почти так же, как обычное расстояние, но отличается от него ослаблением первой аксиомы: из $\text{dist}(f, g) = 0$ не обязательно следует, что $f = g$.² Тогда псевдорасстояние между двумя функциями, совпадающими на заданном наборе значений аргумента, будет равно нулю, даже если эти функции не равны друг другу, т. е. различаются при каких-то других аргументах.

2.2 Интерполирование функций

2.2а Постановка задачи и её свойства

Задача интерполирования — это задача восстановления (доопределения) функции, которая задана на дискретном множестве точек x_i , $i = 0, 1, \dots, n$. Для функций одного вещественного аргумента её формальная постановка такова.

Задан интервал $[a, b] \subset \mathbb{R}$ и конечное множество попарно различных точек $x_i \in [a, b]$, $i = 0, 1, \dots, n$, называемых *узлами интерполяции*. Совокупность всех узлов будем называть *сеткой*. Даны значения y_i , $i = 0, 1, \dots, n$. Требуется построить функцию $g(x)$ от непрерывного аргумента $x \in [a, b]$, которая принадлежит заданному классу функций \mathcal{G} и в узлах x_i принимает значения y_i , $i = 0, 1, \dots, n$. Искомую функцию $g(x)$ называют при этом *интерполирующей функцией* или *интерполянт-ом*.

Практическая значимость задачи интерполяции чрезвычайно велика. Она встречается всюду, где у функции непрерывного аргумента (которая может быть временем, пространственной координатой и т. п.) мы имеем возможность наблюдать лишь значения в дискретном множестве точек, но хотим восстановить по ним ход функции на всём множестве значений аргумента. Например, выполнение многих химических анализов требует существенного времени, так что множество результатов этих анализов по необходимости дискретно. Если нам нужно контролировать по ним непрерывно изменяющийся параметр какого-либо производственного процесса, то неизбежно потребуются интерполирование результатов анализов. Очень часто дискретность множества точек, в

²В отношении этого объекта можно встретить и другие термины. Так, в книге [13] используется термин «квазирасстояние».

которых наблюдаются на практике значения функции, вызвана ограниченностью ресурсов, которые мы можем выделить для сбора данных, или же вообще недоступностью этих данных. Именно так происходит при наблюдении за параметрами земной атмосферы (скоростью и направлением ветра, температурой, влажностью, и пр.) по данным их измерений, которые предоставляются метеостанциями.

В качестве ещё одного примера интерполирования упомянем вычисление различных функций, как элементарных — \sin , \cos , \exp , \log , \dots , так и более сложных, называемых «специальными функциями». С подобной задачей человеческая цивилизация столкнулась очень давно, столетия и даже тысячелетия назад, и типичным способом её решения в докомпьютерную эпоху было составление для нужд практики таблиц — *табуляция* — для значений интересующей нас функции при некоторых специальных фиксированных значениях аргумента, более или менее плотно покрывающих область её определения. Например, \sin и \cos для аргументов, кратных $10'$, т. е. десяти угловым минутам. Подобные таблицы составлялись квалифицированными вычислителями, иногда специально создаваемыми для этой цели организациями, а затем широко распространялись по библиотекам и научным и техническим центрам, где к ним имели доступ люди, занимающиеся практическими вычислениями. Но как, имея подобную таблицу, найти значение интересующей нас функции для аргумента, который не представлен точно в таблице? Скажем, синус угла $17^\circ 23'$ по таблице, где аргумент идёт с шагом $10'$?

Здесь на помощь приходит интерполяция — восстановление значения функции в промежуточных точках по ряду известных значений в некоторых фиксированных опорных точках. Собственно, сам термин «интерполирование» («интерполяция») был впервые употреблён в 1656 году Дж. Валлисом при составлении астрономических и математических таблиц. Он происходит от латинского слова *interpolare*, означающего «переделывать», «подновлять», «ремонттировать».

Для целей практических вычислений таблицы значений различных функций составлялись и издавались вплоть до середины XX века. Вершиной этой деятельности стал выпуск многих томов капитальных таблиц, в которых были тщательно затабулированы все основные функции, встречающиеся в математической и инженерной практике (см., к примеру, [35] и им аналогичные таблицы для других целей).

Интересно, что с появлением и развитием электронных цифровых вычислительных машин описанное применение интерполяции не кану-

ло в лету. В начальный период развития ЭВМ преобладал алгоритмический подход к вычислению элементарных функций, когда основной упор делался на создании алгоритмов, способных «на голом месте» вычислить функцию, исходя из какого-нибудь её аналитического представления, например, в виде быстросходящегося ряда и т. п. (см., к примеру, [19, 20]). Но затем, по мере удешевления памяти ЭВМ и повышения её быстродействия, постепенно распространился подход, очень сильно напоминающий старый добрый табличный способ, но уже на новом уровне. Хранение сотен килобайт или даже мегабайт цифровой информации никаких проблем сейчас не представляет, и потому в современных компьютерах программы вычисления функций (элементарных и специальных), как правило, включают в себя библиотеки за-табулированных значений функции для фиксированных аргументов, опираясь на которые строится значение в нужной нам точке.

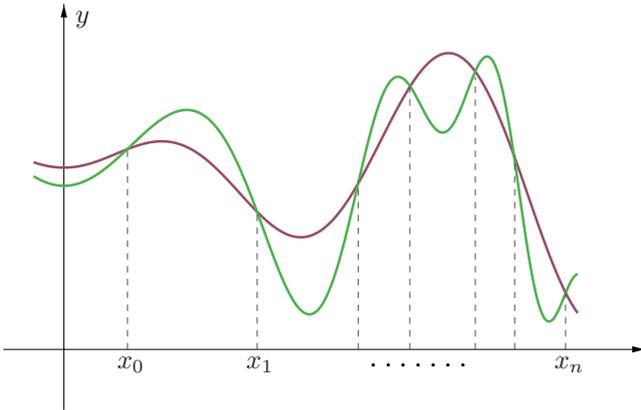


Рис. 2.2. Задача интерполяции может иметь неединственное решение.

Ещё один источник возникновения задачи интерполирования — это желание иметь просто вычисляемое выражение для сложных функциональных зависимостей, которые могут быть заданы как явно, так и неявно, но в исходной форме требуют очень большого труда для своего вычисления.

Если класс \mathcal{G} интерполирующих функций достаточно широк, то решение задачи интерполяции может быть неединственным (см. Рис. 2.2). Напротив, если \mathcal{G} узок, то у задачи интерполяции может вовсе не быть

решений. На практике выбор класса \mathcal{G} обычно диктуется спецификой решаемой практической задачи.

В случае, когда, к примеру, заранее известно, что интерполируемая функция периодична, в качестве интерполирующих функций естественно взять тоже периодические функции — тригонометрические полиномы

$$\sum_{k=0}^m (a_k \cos(kx) + b_k \sin(kx)) \quad (2.4)$$

для некоторого фиксированного m (там, где требуется гладкость), либо пилообразные функции или «ступеньки» (в импульсных системах) и т. п.

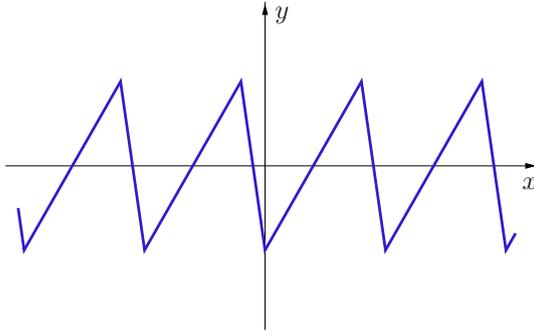


Рис. 2.3. Функция, которую лучше интерполировать с помощью периодических функций.

Ниже мы подробно рассмотрим ситуацию, когда в качестве интерполирующих функций берутся алгебраические полиномы

$$a_0 + a_1x + a_2x^2 + \cdots + a_mx^m, \quad (2.5)$$

которые несложно вычисляются и являются простым и хорошо изученным математическим объектом. При этом мы откладываем до §2.5 рассмотрение вопроса о том, насколько такие полиномы хороши для интерполирования. Вообще, проблема наиболее адекватного выбора класса интерполирующих функций \mathcal{G} не является тривиальной. Для её хорошего решения, как правило, необходимо, чтобы интерполирующие функции были «той же природы», что и интерполируемые функции из

класса \mathcal{F} (который даже не фигурирует в формальной постановке задачи). Если это условие не выполнено, то задача интерполяции может решаться неудовлетворительно.

Определение 2.2.1 Интерполирование функций с помощью алгебраических полиномов называют алгебраической интерполяцией. Алгебраический полином $P_m(x) = a_0 + a_1x + a_2x^2 + \dots + a_mx^m$, решающий задачу алгебраической интерполяции, называется интерполяционным полиномом или алгебраическим интерполянтom.

Как по интерполяционным данным (x_i, y_i) , $i = 0, 1, \dots, n$, найти интерполяционный полином вида (2.5), т.е. определить его коэффициенты?

Подставляя в выражение (2.5) последовательно значения аргумента x_0, x_1, \dots, x_n и учитывая, что получающиеся при этом значения полинома должны быть равны y_0, y_1, \dots, y_n , приходим к соотношениям

$$\begin{array}{rcccl}
 a_0 + a_1x_0 + a_2x_0^2 + \dots + a_mx_0^m & = & y_0, & \\
 a_0 + a_1x_1 + a_2x_1^2 + \dots + a_mx_1^m & = & y_1, & \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\
 a_0 + a_1x_n + a_2x_n^2 + \dots + a_mx_n^m & = & y_n. & & &
 \end{array}
 \tag{2.6}$$

Они образуют систему линейных алгебраических уравнений относительно неизвестных коэффициентов $a_0, a_1, a_2, \dots, a_m$ искомого полинома. Решив её, можно построить и сам полином.

В самом общем случае, если мы не накладываем никаких ограничений на степень полинома m и количество узлов интерполяции $n + 1$, система (2.6) может не иметь решения, а если оно существует, то может быть неединственным. Имеется, тем не менее, важный частный случай задачи алгебраической интерполяции, для которого гарантируется однозначная разрешимость.

Теорема 2.2.1 Если $m = n$, т.е. степень интерполяционного полинома на единицу меньше количества узлов, то решение задачи алгебраической интерполяции существует и единственно.

Доказательство. При $m = n$ матрица системы линейных алгебраиче-

ских уравнений (2.6) — квадратная. Она имеет вид

$$V(x_0, x_1, \dots, x_n) = \begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{pmatrix}, \quad (2.7)$$

и является так называемой матрицей Вандермонда (см., к примеру, [14, 31]). Её определитель равен, как известно, произведению

$$\prod_{1 \leq i < j \leq n} (x_j - x_i),$$

и он не зануляется, если узлы интерполяции попарно отличны друг от друга. Следовательно, система линейных уравнений (2.6) однозначно разрешима тогда при любой правой части, т. е. при любых y_i , $i = 0, 1, \dots, n$. ■

Теорема 2.2.1 и предшествующие ей рассуждения дают, в действительности, конструктивный способ построения интерполяционного полинома через решение системы линейных алгебраических уравнений, который вполне практичен, особенно при небольших n . Он носит общий характер и пригоден для других сходных случаев, когда применяется так называемая *линейная интерполяция*. Этим термином мы будем называть задачу интерполяции, в которой интерполирующие функции из класса \mathcal{G} линейно зависят от некоторых параметров. В частности, это имеет место, когда \mathcal{G} является линейным пространством с заданным базисом. Если же число параметров конечно, то условия удовлетворения интерполяционным данным приводят к необходимости решения системы линейных уравнений, аналогично тому, как это получилось выше.

Если интерполирующие функций из класса \mathcal{G} нельзя представить линейно зависящими от параметров, то соответствующую задачу интерполяции будем называть *нелинейной*. Для определения интерполянта тогда необходимо решать систему нелинейных уравнений.

2.26 Интерполяционный полином Лагранжа

При значительном количестве узлов развитый в предшествующем разделе способ построения интерполянта через решение системы урав-

нений в силу ряда причин невыгоден в вычислительном отношении. Решение систем линейных уравнений само по себе является не вполне тривиальной задачей. Кроме того, система (2.6) оказывается весьма чувствительной к возмущениям данных или, как принято говорить, *плохо обусловленной* (см. §1.3; конкретную числовую оценку чувствительности решения системы (2.6) можно найти в §§3.5а–3.5б). Поэтому получаемый на этом пути интерполяционный полином может обладать большой погрешностью. Наконец, иногда желательно иметь для интерполяционного полинома какое-либо явное аналитическое представление, которого рассмотренный способ всё-таки не даёт.

Систему линейных уравнений (2.6) можно попытаться решить в общем виде с помощью правила Крамера, пользуясь удобным выражением для определителя матрицы Вандермонда в знаменателе и разложением определителей в числителе по столбцу свободных членов $(y_0, y_1, \dots, y_n)^T$. Этот путь может быть успешно пройден, но требует громоздких алгебраических преобразований.

На самом деле нам нечасто требуется знать для интерполяционного полинома коэффициенты канонической формы (2.5). Для большинства практических целей достаточно иметь какое-либо конструктивное представление интерполяционного полинома, позволяющее вычислять его значения в любой наперёд заданной точке.

Для отыскания такого представления заметим, что при фиксированных узлах x_0, x_1, \dots, x_n результат интерполяции линейным образом зависит от значений y_0, y_1, \dots, y_n . Более точно, если полином $P(x)$ решает задачу интерполяции по значениям $y = (y_0, y_1, \dots, y_n)$, а полином $Q(x)$ решает задачу интерполяции с теми же узлами по значениям $z = (z_0, z_1, \dots, z_n)$, то для любых чисел $\alpha, \beta \in \mathbb{R}$ полином $\alpha P(x) + \beta Q(x)$ решает задачу интерполяции для значений $\alpha y + \beta z = (\alpha y_0 + \beta z_0, \alpha y_1 + \beta z_1, \dots, \alpha y_n + \beta z_n)$ на той же совокупности узлов.³

Отмеченным свойством линейности можно воспользоваться для решения задачи интерполяции «по частям», которые удовлетворяют отдельным интерполяционным условиям, а затем собрать эти части воедино. Именно, будем искать интерполяционный полином в виде

$$P_n(x) = \sum_{i=0}^n y_i \phi_i(x), \quad (2.8)$$

³Сказанное можно выразить словами «оператор интерполирования линеен». В действительности, он даже является проектором, и эти наблюдения являются началом большого и плодотворного направления теории приближения функций.

где $\phi_i(x)$ — полином степени n , такой что

$$\phi_i(x_j) = \delta_{ij} = \begin{cases} 0, & \text{при } i \neq j, \\ 1, & \text{при } i = j, \end{cases} \quad (2.9)$$

$i, j = 0, 1, \dots, n$, и посредством δ_{ij} обозначен символ Кронекера. Полином $y_i \phi_i(x)$ имеет степень n и решает задачу интерполяции для значений $(0, \dots, 0, y_i, 0, \dots, 0)$ в узлах $x_0, x_1, \dots, x_n, i = 0, 1, \dots, n$, и потому полином $P_n(x)$, задаваемый представлением (2.8), в целом действительно удовлетворяет условиям задачи.

Коль скоро $\phi_i(x)$ зануляется в точках $x_0, \dots, x_{i-1}, x_{i+1}, \dots, x_n$, то ясно, что он должен иметь вид

$$\phi_i(x) = \Phi_i(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n). \quad (2.10)$$

В правой части этого равенства произведение n линейных по x членов даёт полином степени n , так что Φ_i должно быть некоторым числовым множителем. Для его определения подставим в выражение (2.10) значение аргумента $x = x_i$, откуда в силу (2.9) получается

$$\Phi_i(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n) = 1.$$

Следовательно,

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

и потому

$$\phi_i(x) = \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)}.$$

Полиномы $\phi_i(x)$ называют *базисными полиномами Лагранжа*, а иногда также *полиномами влияния i -го узла* (последний термин объясняется условием (2.9)).

В целом, из (2.8) следует, что задачу алгебраической интерполяции решает полином

$$P_n(x) = \sum_{i=0}^n y_i \frac{\prod_{j \neq i} (x - x_j)}{\prod_{j \neq i} (x_i - x_j)}. \quad (2.11)$$

Его называют *интерполяционным полиномом в форме Лагранжа* или просто *интерполяционным полиномом Лагранжа*.

Далее нам потребуется его запись в несколько другом виде. Введём вспомогательную функцию

$$\omega_n(x) = (x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n) \quad (2.12)$$

— полином $(n + 1)$ -й степени, зануляющийся во всех узлах интерполяции. Тогда

$$\phi_i(x) = \frac{\omega_n(x)}{(x - x_i) \omega'_n(x_i)}, \quad (2.13)$$

и поэтому

$$P_n(x) = \sum_{i=0}^n y_i \frac{\omega_n(x)}{(x - x_i) \omega'_n(x_i)}. \quad (2.14)$$

Задача интерполяции полностью решается с помощью полиномов (2.11) и (2.14), которые находят широчайшее применение в вычислительной практике. Тем не менее, в ряде случаев они оказываются не совсем удобными. Дело в том, что каждый из базисных полиномов Лагранжа $\phi_i(x)$ зависит от всех узлов интерполяции сразу. По этой причине если, к примеру, мы имеем дело с изменяющимся набором узлов, то каждый раз должны будем перевычислять все $\phi_i(x)$. Иными словами, при смене набора узлов интерполяции полином Лагранжа претерпевает большое изменение и должен быть перевычислен заново.

Нельзя ли найти такую форму интерполяционного полинома, которая изменялась бы незначительно при небольших изменениях в наборе узлов интерполяции? Этот вопрос решается с помощью интерполяционного полинома в форме Ньютона, и для его построения нам будет необходима новая техника, основанная на понятии разделённой разности от функции.

2.2в Разделённые разности и их свойства

Пусть дана функция f и попарно различные точки x_0, x_1, \dots, x_n из её области определения, в которых функция принимает значения $f(x_0), f(x_1), \dots, f(x_n)$. *Разделёнными разностями* функции f , обозначаемыми $f^{\angle}(x_i, x_{i+1})$, называются отношения

$$f^{\angle}(x_i, x_{i+1}) := \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}, \quad (2.15)$$

$i = 0, 1, \dots, n - 1$. Их называют также разделёнными разностями первого порядка.

Разделённые разности второго порядка — это величины

$$f^{\angle}(x_i, x_{i+1}, x_{i+2}) := \frac{f^{\angle}(x_{i+1}, x_{i+2}) - f^{\angle}(x_i, x_{i+1})}{x_{i+2} - x_i}, \quad (2.16)$$

$i = 0, 1, \dots, n - 2$, которые являются разделёнными разностями от разделённых разностей. Аналогичным образом вводятся разделённые разности высших порядков: *разделённая разность k -го порядка* от функции f есть, по определению,

$$f^{\angle}(x_i, x_{i+1}, \dots, x_{i+k}) := \frac{f^{\angle}(x_{i+1}, \dots, x_{i+k}) - f^{\angle}(x_i, \dots, x_{i+k-1})}{x_{i+k} - x_i}, \quad (2.17)$$

$i = 0, 1, \dots, n - k$, т.е. она равна разделённой разности от разделённых разностей предыдущего $(k - 1)$ -го порядка. Для удобства и единообразия можно также считать, что сами значения функции являются разделёнными разностями нулевого порядка, т.е. $f^{\angle}(x_i) = f(x_i)$, $i = 0, 1, \dots, n$.

Разделённые разности можно определять не только для функций непрерывного аргумента, но и для функций дискретного аргумента, проще говоря, для набора значений y_0, y_1, \dots, y_n , соответствующего узлам x_0, x_1, \dots, x_n . Назовём разделённой разностью первого порядка между узлами x_i и x_{i+1} величину

$$(y_i, y_{i+1})^{\angle} := \frac{y_{i+1} - y_i}{x_{i+1} - x_i}.$$

Разделённой разностью k -го порядка значений $y_i, y_{i+1}, \dots, y_{i+k}$ по узлам $x_i, x_{i+1}, \dots, x_{i+k}$ называется величина

$$(y_i, y_{i+1}, \dots, y_{i+k})^{\angle} := \frac{(y_{i+1}, \dots, y_{i+k})^{\angle} - (y_i, \dots, y_{i+k-1})^{\angle}}{x_{i+k} - x_i},$$

$i = 0, 1, \dots, n - k$. Это обозначение не содержит явного указания на узлы $x_i, x_{i+1}, \dots, x_{i+k}$, по которым берётся набор значений $(y_i, y_{i+1}, \dots, y_{i+k})$, так что значения этих узлов подразумеваются.

Отметим, что в определении разделённых разностей, вообще говоря, не накладываются никаких условий на взаимное расположение точек x_0, x_1, \dots, x_n . В частности, совсем не обязательно, чтобы $x_i < x_{i+1}$.

Понятию разделённой разности можно также придать смысл для случая совпадающих узлов $x_i = x_{i+1}$, если понимать его как результат предельного перехода при $x_i \rightarrow x_{i+1}$ (см. подробности, к примеру, в [17, 56]).

Нетрудно увидеть геометрический смысл разделённой разности первого порядка. Будучи отношением приращения функции к приращению её аргумента, это угловой коэффициент (тангенс угла наклона к оси абсцисс) секущей графика функции $y = f(x)$, взятой между точками с аргументами x_i и x_{i+1} . В общем случае разделённая разность функции — это «средняя скорость» её изменения на рассматриваемом интервале, в отличие от «мгновенной скорости» изменения функции в точке, выражаемой её производной $f'(x)$.

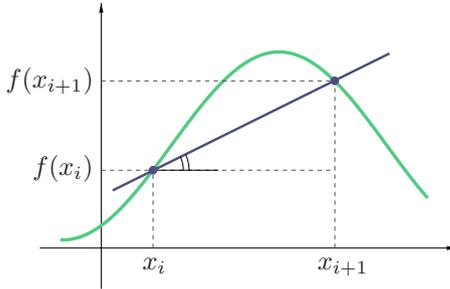


Рис. 2.4. Иллюстрация смысла разделённых разностей, как углового коэффициента секущей графика функции

Если \tilde{x} — какая-то фиксированная точка, то для любой другой точки x имеет место равенство

$$f(x) = f(\tilde{x}) + f'(\tilde{x}, x)(x - \tilde{x}),$$

аналогичное формуле Тейлора, в которой удержаны лишь члены первого порядка. В связи с этим уместно упомянуть, что разделённую разность иногда называют *наклоном* функции между заданными точками (см. [13]). Разделённые разности-наклоны могут быть определены для функций многих переменных и даже для операторов, действующих из одного абстрактного пространства в другое. Интересно, что в начале XX века для обозначения этой конструкции использовался также термин «подъём функции» [64]. В математических текстах для разделённых разностей функции f по точкам $x_i, x_{i+1}, \dots, x_{i+k}$ нередко

используется обозначение $f[x_i, x_{i+1}, \dots, x_{i+k}]$ или даже маловыразительное $f(x_i, x_{i+1}, \dots, x_{i+k})$.

Операция взятия разделённой разности является линейной: для любых функций f, g и для любых скаляров α, β справедливо

$$(\alpha f + \beta g)^\zeta = \alpha f^\zeta + \beta g^\zeta \quad (2.18)$$

при одинаковых аргументах разделённых разностей. Это очевидно следует из определения для разделённой разности первого порядка, а для разделённых разностей высших порядков показывается несложной индукцией по величине порядка. То же самое верно и для разделённых разностей от наборов значений:

$$(\alpha(y_i, \dots, y_{i+k}) + \beta(z_i, \dots, z_{i+k}))^\zeta = \alpha(y_i, \dots, y_{i+k})^\zeta + \beta(z_i, \dots, z_{i+k})^\zeta$$

по одному и тому же набору узлов.

Предложение 2.2.1 *Имеет место представление*

$$f^\zeta(x_i, x_{i+1}, \dots, x_{i+k}) = \sum_{j=i}^{i+k} \frac{f(x_j)}{\prod_{\substack{l=i \\ l \neq j}}^{i+k} (x_j - x_l)}. \quad (2.19)$$

Для разделённой разности от набора значений аналогичная формула выглядит следующим образом

$$(y_i, y_{i+1}, \dots, y_{i+k})^\zeta = \sum_{j=i}^{i+k} \frac{y_j}{\prod_{\substack{l=i \\ l \neq j}}^{i+k} (x_j - x_l)}.$$

Доказательство. Оно проводится индукцией по порядку k разделённой разности.

При $k = 1$ доказываемая формула, как нетрудно проверить, совпадает с определением разделённой разности первого порядка.

Пусть Предложение уже доказано для некоторого положительного

целого k . Тогда для разделённой разности $k+1$ -го порядка будем иметь

$$f^{\prime\prime}(x_i, x_{i+1}, \dots, x_{i+k+1})$$

$$= \frac{f^{\prime\prime}(x_{i+1}, x_{i+1}, \dots, x_{i+k+1}) - f^{\prime\prime}(x_i, x_{i+1}, \dots, x_{i+k})}{x_{i+k+1} - x_i}$$

по определению разделённой разности

$$= \frac{1}{x_{i+k+1} - x_i} \cdot \left(\sum_{j=i+1}^{i+k+1} \frac{f(x_j)}{\prod_{\substack{l=i+1 \\ l \neq j}}^{i+k+1} (x_j - x_l)} - \sum_{j=i}^{i+k} \frac{f(x_j)}{\prod_{\substack{l=i \\ l \neq j}}^{i+k} (x_j - x_l)} \right)$$

согласно индукционному предположению

$$= \frac{f(x_{i+k+1})}{(x_{i+k+1} - x_i) \prod_{l=i+1}^{i+k} (x_{i+k+1} - x_l)}$$

$$+ \frac{1}{x_{i+k+1} - x_i} \cdot \left(\sum_{j=i+1}^{i+k} \frac{f(x_j)}{\prod_{\substack{l=i+1 \\ l \neq j}}^{i+k+1} (x_j - x_l)} - \sum_{j=i+1}^{i+k} \frac{f(x_j)}{\prod_{\substack{l=i \\ l \neq j}}^{i+k} (x_j - x_l)} \right)$$

$$- \frac{f(x_i)}{(x_{i+k+1} - x_i) \prod_{l=i+1}^{i+k} (x_i - x_l)}$$

после вынесения из-под скобок первого слагаемого первой суммы и последнего слагаемого второй суммы. В полученную сумму члены с $f(x_i)$ и $f(x_{i+k+1})$ — первый и последний — входят по одному разу, причём их коэффициенты уже имеют тот вид, который утверждается в Предложении. Члены с остальными $f(x_j)$ входят дважды, и после приведения

подобных членов коэффициент при $f(x_j)$ будет равен

$$\begin{aligned} & \frac{1}{x_{i+k+1} - x_i} \cdot \left(\frac{1}{\prod_{\substack{l=i+1 \\ l \neq j}}^{i+k+1} (x_j - x_l)} - \frac{1}{\prod_{\substack{l=i \\ l \neq j}}^{i+k} (x_j - x_l)} \right) \\ &= \frac{(x_j - x_i) - (x_j - x_{i+k+1})}{(x_{i+k+1} - x_i) \prod_{\substack{l=i \\ l \neq j}}^{i+k+1} (x_j - x_l)} = \frac{1}{\prod_{\substack{l=i \\ l \neq j}}^{i+k+1} (x_j - x_l)}, \end{aligned}$$

что и требовалось показать. ■

Следствие. Разделённая разность — симметричная функция своих аргументов, т. е. узлов x_0, x_1, \dots, x_k . Иными словами, она не изменяется при любой их перестановке. Это непосредственно следует из симметричного вида выражения, стоящего в правой части (2.19).

Для фиксированного набора узлов численные значения разделённых разностей любой функции нетрудно вычислить согласно определениям (2.15)–(2.16) или по формуле (2.19). Но сложность выражений для разделённых разностей как функций узлов в общем случае быстро возрастает с ростом порядка разделённой разности. Тем не менее, в случае алгебраических полиномов выражения для разделённых разностей относительно просто получаются из выражений для исходной функции. Вспомним известную формулу элементарной алгебры $x^n - y^n = (x - y)(x^{n-1} + x^{n-2}y + \dots + xy^{n-2} + y^{n-1})$, из которой следует, что

$$\frac{x^n - y^n}{x - y} = x^{n-1} + x^{n-2}y + \dots + xy^{n-2} + y^{n-1}. \quad (2.20)$$

Этот результат позволяет явно выписать разделённую разность для любой целой степени переменной. Далее для произвольного полинома можно воспользоваться свойством линейности разделённой разности (2.18).

Пример 2.2.1 Вычислим разделённые разности от полинома $g(x) = x^3 - 4x + 1$.

Будем искать по отдельности разделённые разности от мономов, образующих $g(x)$. В силу (2.20) имеем

$$\frac{x_2^3 - x_1^3}{x_2 - x_1} = x_2^2 + x_2x_1 + x_1^2.$$

Для линейного монома $(-4x)$ разделённая разность находится тривиально и равна (-4) , а для константы 1 она равна нулю. Следовательно, в целом

$$g'(x_1, x_2) = x_2^2 + x_2x_1 + x_1^2 - 4.$$

Вычислим вторую разделённую разность от $g(x)$:

$$\begin{aligned} g''(x_1, x_2, x_3) &= \frac{g'(x_2, x_3) - g'(x_1, x_2)}{x_3 - x_1} \\ &= \frac{(x_3^2 + x_3x_2 + x_2^2 - 4) - (x_2^2 + x_2x_1 + x_1^2 - 4)}{x_3 - x_1} \\ &= \frac{x_3^2 + (x_3 - x_1)x_2 - x_1^2}{x_3 - x_1} = x_1 + x_2 + x_3. \end{aligned}$$

Третья разделённая разность

$$\begin{aligned} g'''(x_1, x_2, x_3, x_4) &= \frac{g''(x_2, x_3, x_4) - g''(x_1, x_2, x_3)}{x_4 - x_1} \\ &= \frac{(x_2 + x_3 + x_4) - (x_1 + x_2 + x_3)}{x_4 - x_1} \\ &= \frac{x_4 - x_1}{x_4 - x_1} = 1, \end{aligned}$$

т. е. является постоянной. Четвёртая и последующие разделённые разности от $g(x)$ будут, очевидно, тождественно нулевыми функциями. ■

Как видим, взятие разделённой разности от полинома уменьшает его степень на единицу, так что разделённые разности порядка более n от полинома степени n равны нулю. Сделанное наблюдение демонстрирует глубокую аналогию между разделёнными разностями и производными: каждое применение дифференцирования к полиному также последовательно уменьшает его степень на единицу. В действительности, эта связь видна даже из определения разделённой разности первого

порядка, которую можно рассматривать как «недоделанную производную», поскольку у неё отсутствует предельный переход одного аргумента к другому.

Предложение 2.2.2 (связь разделённых разностей с производными)
Пусть $f \in C^n[a, b]$, т. е. функция f непрерывно дифференцируема n раз на интервале $[a, b]$, где расположены узлы x_0, x_1, \dots, x_n , и пусть $\underline{x} = \min\{x_0, x_1, \dots, x_n\}$, $\bar{x} = \max\{x_0, x_1, \dots, x_n\}$. Тогда

$$f^{\zeta}(x_0, x_1, \dots, x_n) = \frac{1}{n!} f^{(n)}(\xi)$$

для некоторой точки $\xi \in]\underline{x}, \bar{x}[$.

Для разделённых разностей первого порядка этот факт непосредственно следует из теоремы Лагранжа о среднем (формулы конечных приращений), согласно которой

$$f(x_{i+1}) - f(x_i) = f'(\xi) \cdot (x_{i+1} - x_i)$$

для некоторой точки $\xi \in]x_i, x_{i+1}[$. Для общего случая доказательство Предложения 2.2.2 будет приведено несколько позже, в §2.2д.

Существует более точное (хотя и более громоздкое) интегральное представление для разделённых разностей, о котором можно подробно узнать в [17, 53, 66].

2.2г Интерполяционный полином Ньютона

Выведем теперь другую форму интерполяционного полинома с учётом требования иметь такое выражение, которое в минимальной степени перестраивалось бы при смене набора узлов интерполяции.

Обозначим через $P_k(x)$ интерполяционный полином степени k , построенный по узлам x_0, x_1, \dots, x_k . В частности, $P_0(x) = f(x_0)$, имея нулевую степень и будучи построенным по одному узлу x_0 . Тогда очевидно следующее тождество

$$P_n(x) = P_0(x) + \sum_{k=1}^n (P_k(x) - P_{k-1}(x)). \quad (2.21)$$

Замечательность этого представления состоит в том, что при добавлении или удалении последних по номеру узлов интерполяции перестройке должны подвергнуться лишь те последние слагаемые суммы

из правой части (2.21), которые вовлекает эти изменяемые узлы. Первые слагаемые в (2.21) зависят только от первых узлов интерполяции и останутся неизменными, Таким образом, стоящая перед нами задача окажется решённой, если будут найдены удобные и просто выписываемые выражения для разностей $P_k(x) - P_{k-1}(x)$.

Заметим, что разность $(P_k(x) - P_{k-1}(x))$ есть полином степени k , который обращается в нуль в узлах x_0, x_1, \dots, x_{k-1} , общих для $P_k(x)$ и $P_{k-1}(x)$, где эти полиномы должны принимать одинаковые значения $f(x_0), f(x_1), \dots, f(x_{k-1})$. Поэтому должно быть

$$P_k(x) - P_{k-1}(x) = A_k(x - x_0)(x - x_1) \cdots (x - x_{k-1})$$

с некоторой константой A_k . Для её определения вспомним, что по условию интерполяции $P_k(x_k) = y_k$. Следовательно,

$$A_k = \frac{y_k - P_{k-1}(x_k)}{(x_k - x_0)(x_k - x_1) \cdots (x_k - x_{k-1})}.$$

Отсюда, подставляя вместо $P_{k-1}(x)$ выражение для интерполяционного полинома в форме Лагранжа, нетрудно вывести, что

$$\begin{aligned} A_k &= \frac{1}{\prod_{l=0}^{k-1} (x_k - x_l)} \cdot \left(y_k - \sum_{j=0}^{k-1} y_j \frac{\prod_{\substack{l=0 \\ l \neq j}}^{k-1} (x_k - x_l)}{\prod_{\substack{l=0 \\ l \neq j}}^{k-1} (x_j - x_l)} \right) \\ &= \frac{y_k}{\prod_{l=0}^{k-1} (x_k - x_l)} - \sum_{j=0}^{k-1} \left(\frac{1}{\prod_{l=0}^{k-1} (x_k - x_l)} y_j \frac{\prod_{\substack{l=0 \\ l \neq j}}^{k-1} (x_k - x_l)}{\prod_{\substack{l=0 \\ l \neq j}}^{k-1} (x_j - x_l)} \right) \\ &= \frac{y_k}{\prod_{l=0}^{k-1} (x_k - x_l)} - \sum_{j=0}^{k-1} \frac{y_j}{(x_k - x_j) \prod_{\substack{l=0 \\ l \neq j}}^{k-1} (x_j - x_l)} \quad \begin{array}{l} \text{после сокращения} \\ \text{произведений} \end{array} \end{aligned}$$

$$\begin{aligned}
&= \frac{y_k}{\prod_{\substack{l=0 \\ l \neq k}}^k (x_k - x_l)} + \sum_{j=0}^{k-1} \frac{y_j}{(x_j - x_k) \prod_{\substack{l=0 \\ l \neq j}}^{k-1} (x_j - x_l)} \\
&= \sum_{j=0}^k \frac{y_j}{\prod_{\substack{l=0 \\ l \neq j}}^k (x_j - x_l)} = (y_0, y_1, \dots, y_k)^\zeta
\end{aligned}$$

в силу представления, доказанного в Предложении 2.2.1. Окончательно

$$\begin{aligned}
P_n(x) = & y_0 + (y_0, y_1)^\zeta (x - x_0) + (y_0, y_1, y_2)^\zeta (x - x_0)(x - x_1) + \\
& \dots + (y_0, y_1, y_2, \dots, y_n)^\zeta (x - x_0)(x - x_1) \dots (x - x_{n-1}).
\end{aligned}$$

Выражение в правой части этого равенства называется интерполяционным полиномом в форме Ньютона, или просто *интерполяционным полиномом Ньютона*. Оно является равносильной формой записи интерполяционного полинома, широко применяемой в ситуациях, где использование формы Лагранжа по тем или иным причинам оказывается неудобным. Полезно иметь в виду следующее представление⁴

$$\begin{aligned}
P_n(x) = P_k(x) + (y_0, y_1, \dots, y_{k+1})^\zeta (x - x_0) \dots (x - x_k) \\
+ \dots + (y_0, y_1, \dots, y_n)^\zeta (x - x_0) \dots (x - x_{n-1}),
\end{aligned} \tag{2.22}$$

справедливое для любого k , такого что $0 \leq k \leq n - 1$.

Пусть f — вещественная n -гладкая функция. С учётом результата Предложения 2.2.2, т. е. равенства

$$f^\zeta(x_0, x_1, \dots, x_n) = \frac{1}{n!} f^{(n)}(\xi),$$

хорошо видно, что интерполяционный полином Ньютона для гладкой функции непрерывного аргумента является прямым аналогом формулы Тейлора (полинома Тейлора)

$$f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!} (x - x_0)^2 + \dots + \frac{f^{(n)}(x_0)}{n!} (x - x_0)^n.$$

⁴Образно выражаясь, оно показывает, как интерполяционные полиномы Ньютона разных степеней вложены друг в друга наподобие «матрёшек».

При этом аналогами степеней переменной $(x - x_0)^k$ являются произведения $(x - x_0)(x - x_1) \cdots (x - x_k)$, которые в случае равномерно расположенных и упорядоченных по возрастанию узлов x_0, x_1, \dots, x_k часто называют *обобщённой степенью* [9].

Практическое нахождение интерполяционного полинома Ньютона требует знания всех разделённых разностей функции, и наиболее удобно вычислять их всё-таки по рекуррентным формулам (2.15)–(2.17).

Важнейший частный случай интерполирования относится к равномерному расположению узлов, когда величина $h_i = x_i - x_{i+1}$ постоянна и не зависит от i . Тогда вычисление разделённых разностей решительно упрощается, сводясь к оперированию с так называемыми *конечными разностями*. По определению конечной разностью (иногда добавляют — первого порядка) от функции f в точке x называется величина

$$\Delta y = \Delta f(x) = f(x + h) - f(x).$$

В частности, можно считать, что $\Delta x = h$. Конечные разности второго порядка $\Delta^2 f(x)$ — это конечные разности от конечных разностей, и далее рекуррентно.

Таблица 2.1. Горизонтальная таблица конечных разностей функции

x	y	Δy	$\Delta^2 y$	\dots	$\Delta^n y$
x_0	y_0	Δy_0	$\Delta^2 y_0$	\dots	$\Delta^n y_0$
x_1	y_1	Δy_1	$\Delta^2 y_1$	\dots	$\Delta^n y_1$
x_2	y_2	Δy_2	$\Delta^2 y_2$	\dots	$\Delta^n y_2$
\vdots	\vdots	\vdots	\vdots	\ddots	\vdots

Интерполяционный полином Ньютона для равномерно расположен-

ных узлов принимает вид

$$P_n(x) = f(x_0) + \frac{\Delta f(x_0)}{1!h} (x - x_0) + \frac{\Delta^2 f(x_0)}{2!h^2} (x - x_0)(x - x_1) + \dots + \frac{\Delta^n f(x_0)}{n!h^n} (x - x_0)(x - x_1) \cdots (x - x_{n-1}),$$

особенно сильно похожий на полином Тейлора, тем более, что классическое обозначение для производных $f^{(n)} = \frac{d^n f}{dx^n}$.

Вычисление конечных разностей таблично заданной функции удобно оформлять также в виде таблицы (см. Табл. 2.1), где в дополнительных столбцах, заполняемых последовательно один за другим (слева направо), выписываются значения конечных разностей.

2.2д Погрешность алгебраической интерполяции с простыми узлами

Задача интерполяции, успешно решённая в предшествующих параграфах, часто находится в более широком контексте, описанном во введении к этой теме, на стр. 41. Именно, значения y_0, y_1, \dots, y_n принимаются в узлах x_0, x_1, \dots, x_n некоторой реальной функцией непрерывного аргумента $f(x)$, свойства которой (хотя бы отчасти) известны. Насколько сильно будет отличаться от неё построенный нами интерполянт? Именно это отличие понимается под «погрешностью интерполяции».

Определение 2.2.2 *Остаточным членом или остатком интерполяции называется функция $R(f, x) = f(x) - g(x)$, являющаяся разностью рассматриваемой функции $f(x)$ и интерполирующей её функции $g(x)$.*

Предложение 2.2.3 *Если точка z не совпадает ни с одним из узлов интерполирования x_0, x_1, \dots, x_n , то в задаче алгебраической интерполяции остаточный член $R_n(f, z) := f(z) - P_n(z)$ равен*

$$R_n(f, z) = f^{(n)}(x_0, x_1, \dots, x_n, z) \cdot \omega_n(z), \quad (2.23)$$

где функция ω_n определяется посредством (2.12), т. е.

$$\omega_n(x) = \prod_{i=0}^n (x - x_i).$$

Доказательство. Выпишем для f интерполяционный полином Ньютона $(n + 1)$ -й степени по узлам x_0, x_1, \dots, x_n, z . Согласно (2.22)

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

где $P_n(x)$ — полином Ньютона для узлов x_0, x_1, \dots, x_n . Подставляя в это соотношение значение $x = z$, получим

$$P_{n+1}(z) = P_n(z) + f^{\zeta}(x_0, x_1, \dots, x_n, z)(z - x_0)(z - x_1) \cdots (z - x_n).$$

Но $P_{n+1}(z) = f(z)$ по построению полинома P_{n+1} . Поэтому

$$\begin{aligned} R_n(f, z) &= f(z) - P_n(z) \\ &= f^{\zeta}(x_0, x_1, \dots, x_n, z)(z - x_0)(z - x_1) \cdots (z - x_n), \end{aligned}$$

что и требовалось. ■

Полученный результат позволяет точно находить численное значение погрешности интерполирования в конкретных точках, но он не очень полезен для исследования поведения погрешности «в целом», на всём интервале интерполирования. Чтобы получить более удобные оценки для остаточного члена, нам будет необходимо воспользоваться Предложением 2.2.2 о связи разделённых разностей и производных, и мы дадим его строгое доказательство.

Доказательство Предложения 2.2.2.

Поскольку разделённая разность есть симметричная функция узлов, то в формуле (2.2.2) без какого-либо ограничения общности можно считать эти узлы x_0, x_1, \dots, x_n упорядоченными по возрастанию индекса, т. е. $x_0 < x_1 < \dots < x_n$. Обозначив

$$\theta(x) := f^{(n)}(x) - n! f^{\zeta}(x_0, x_1, \dots, x_n),$$

можно заметить, что Предложение 2.2.2 равносильно тогда следующему утверждению: на] x_0, x_n [существует точка ξ , которая является нулём функции $\theta(x)$.

По точкам x_0, x_1, \dots, x_n построим для функции $f(x)$ интерполяционный полином $P_n(x)$. Тогда введённая выше функция $\theta(x)$ есть n -ая производная по x от остаточного члена интерполяции $R_n(f, x) = f(x) - P_n(x)$, т. е.

$$f^{(n)}(x) - n! f^{\zeta}(x_0, x_1, \dots, x_n) = R_n^{(n)}(f, x),$$

в чём можно убедиться непосредственным дифференцированием равенства

$$R_n(f, x) = f(x) - P_n(x),$$

где интерполяционный полином $P_n(x)$ берётся в форме Ньютона.

В самом деле, в интерполяционном полиноме Ньютона только у разделённой разности n -го порядка $f^{\angle}(x_0, x_1, \dots, x_n)$ коэффициент является полиномом n -ой степени со старшим членом x^n . Коэффициенты остальных разделённых разностей — полиномы меньших степеней от x , которые исчезнут при n -кратном дифференцировании, тогда как от полинома n -ой степени со старшим членом x^n после этого дифференцирования останется число $n!$.

Функция $R_n(f, x)$ является n -кратно дифференцируемой на $[a, b]$ и, кроме того, обращается в нуль в $n+1$ различных точках — узлах интерполяции x_0, x_1, \dots, x_n . В силу известной из математического анализа теоремы Ролля производная $R'_n(f, x)$ обязана зануляться внутри n интервалов $[x_0, x_1], [x_1, x_2], \dots, [x_{n-1}, x_n]$, т. е. она имеет n нулей.

Далее, повторяя те же рассуждения в отношении второй производной $R''_n(f, x)$, приходим к выводу, что она должна иметь на $]x_0, x_n[$ не менее $n-1$ нулей. Аналогично для третьей производной $R'''_n(f, x)$ и т. д. вплоть до $R^{(n)}(f, x)$, которая должна иметь на $]x_0, x_n[$ хотя бы один нуль. Это и требовалось доказать. ■

Теорема 2.2.2 Пусть $f \in C^{n+1}[a, b]$, т. е. функция $f(x)$ непрерывно дифференцируема $n+1$ раз на интервале $[a, b]$. При её интерполировании по попарно различным узлам x_0, x_1, \dots, x_n с помощью полинома n -ой степени остаточный член $R_n(f, x)$ может быть представлен в виде

$$R_n(f, x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \cdot \omega_n(x), \quad (2.24)$$

где $\xi(x)$ — некоторая точка, принадлежащая открытому интервалу $]a, b[$ и зависящая от x , а $\omega_n = (x-x_0)(x-x_1)\dots(x-x_n)$.

Доказательство. Если $x = x_i$ для одного из узлов интерполирования, то $R_n(f, x) = 0$, но в то же время и $\omega_n(x) = 0$. Поэтому в качестве ξ в этом случае можно взять любую точку из открытого интервала $]a, b[$.

Если же аргумент x остаточного члена не совпадает ни с одним из узлов интерполирования, то применяем Предложение 2.2.3, в котором разделённая разность представлена согласно результату Предложения 2.2.2. ■

Выражение (2.24) для остаточного члена алгебраической интерполяции обычно связывают с именем О.Л. Коши, впервые его получившего. Другое выражение для остаточного члена, не использующее неизвестную точку $\xi(x)$ и основанное на интегральном представлении разделённых разностей, можно найти, к примеру, в книгах [17, 66].

Если обозначить

$$M_n = \max_{\xi \in [a, b]} |f^{(n)}(\xi)|$$

— максимум абсолютного значения n -ой производной на рассматриваемом интервале, то нетрудно выписать огрублённые оценки, вытекающие из (2.24) и полезные при практическом вычислении погрешности интерполирования:

$$|R_n(f, x)| \leq \frac{M_{n+1}}{(n+1)!} \cdot |\omega_n(x)|, \quad (2.25)$$

или даже совсем простую

$$|R_n(f, x)| \leq \frac{M_{n+1}(b-a)^{n+1}}{(n+1)!}. \quad (2.26)$$

Для оценивания максимума $(n+1)$ -ой производной функции можно воспользоваться, к примеру, интервальными методами, взяв какое-либо интервальное расширение для $f^{(n+1)}(x)$ на $[a, b]$ (см. §1.5).

Отметим, что полученные выше оценки — (2.24) и её следствия (2.25) и (2.26) — становятся неприменимыми, если функция f имеет гладкость, меньшую чем $n+1$. В то же время представление погрешности интерполирования в виде (2.23) работает для любых функций.

В представлении (2.24) поведение полинома $\omega_n(x)$ при изменении x типично для полиномов с вещественными корнями вообще. Пусть, как и ранее, $\underline{x} = \min\{x_0, x_1, \dots, x_n\}$, $\bar{x} = \max\{x_0, x_1, \dots, x_n\}$. Если аргумент x находится на интервале $[\underline{x}, \bar{x}]$ расположения корней x_0, x_1, \dots, x_n или «не слишком далёко» от него, то $\omega_n(x)$ принимает относительно умеренные значения, так как формирующие его множители $(x - x_i)$, $i = 0, 1, \dots, n$, «не слишком сильно» отличаются от нуля. Если же значения аргумента x находятся на существенном удалении от корней полинома $\omega_n(x)$, то его абсолютная величина и вместе с ней погрешность алгебраической интерполяции, очень быстро растут. На Рис. 2.5 изображён пример графика такого полинома нечётной (седьмой) степени.

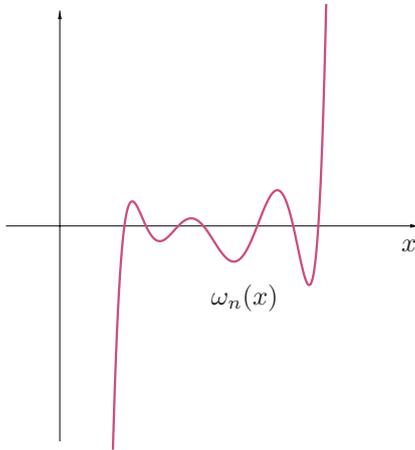


Рис. 2.5. Типичное поведение полинома $\omega_n(x)$: быстрый рост за пределами интервала узлов

В связи со сказанным выше полезно на качественном уровне различать два случая. Если значения интерполируемой функции ищутся в точках, далёких от интервала узлов интерполяции, используют термин *экстраполяция*. Ей противопоставляется *интерполяция* в узком смысле, когда значения функции восстанавливаются на интервале, где расположены узлы, или же вблизи от него. Из наших рассуждений следует, что экстраполяция, как правило, сопровождается существенными ошибками, и потому не стоит использовать её слишком широко.

В рассмотренной выше постановке задачи интерполирования (§2.2а) расположение узлов считалось данным извне и фиксированным. Подобный подход соответствует тем практическим задачам, в которых измерения величины y_i могут осуществляться, к примеру, лишь в какие-то фиксированные моменты времени x_i , либо в определённых выделенных точках пространства и т. п., то есть заданы каким-то внешним образом и не могут быть изменены по нашему желанию.

Но существуют задачи, в которых мы можем управлять выбором узлов интерполирования. При этом естественно возникает вопрос о том, как сделать этот выбор наилучшим образом, чтобы погрешность интерполирования была как можно меньшей. В наиболее общей формулировке эта задача является весьма трудной, и её решение существенно

завязано на свойства интерполируемой функции $f(x)$. Но имеет смысл рассмотреть и упрощённую постановку задачи, в которой на заданном интервале минимизируются значения полинома $\omega_n(x)$, тогда как множитель $f^z(x_0, x_1, \dots, x_n, z)$ или множитель $f^{(n+1)}(\xi(x))/(n+1)!$ в выражениях для остаточного члена (2.23) или (2.24) соответственно округлённо считаются «приближёнными константами».

Фактически, ответ на поставленный вопрос сводится к подбору узлов x_0, x_1, \dots, x_n в пределах заданного интервала $[a, b]$ так, чтобы полином $\omega_n(x) = (x - x_0)(x - x_1) \dots (x - x_n)$ принимал «как можно меньшие значения» на $[a, b]$. Конкретный смысл, который вкладывается в эту фразу, может быть весьма различен, так как функция — полином $\omega_n(x)$ в нашем случае — определяется своими значениями в бесконечном множестве аргументов, и малость одних значений функции может иметь место наряду с очень большими значениями при других аргументах (см. Рис. 2.18). Ниже мы рассматриваем ситуацию, когда «отклонение от нуля» понимается как равномерное (чебышёвское) расстояние (2.1) до нулевой функции, т. е. как максимум абсолютных значений функции на интервале. Это условие является одним из наиболее часто встречающихся в прикладных задачах.

2.3 Полиномы Чебышёва

2.3а Определение и основные свойства

Полиномы Чебышёва — это семейство полиномов, обозначаемых по традиции⁵ как $T_n(x)$, и зависящих от неотрицательного целого параметра n . Они могут быть определены различными равносильными способами, и наиболее просто и наглядно их *тригонометрическое определение*:

$$T_n(x) = \cos(n \arccos x), \quad x \in [-1, 1], \quad (2.27)$$

$n = 0, 1, 2, \dots$. Как известно, всякий полином степени n однозначно определяется своими значениями в $(n+1)$ точках, а формулой (2.27) мы фактически задаём значения функции в бесконечном множестве точек из $[-1, 1]$. Поэтому если посредством (2.27) на $[-1, 1]$ в самом деле задаются полиномы, то они однозначно определяются с помощью

⁵С буквы «Т» начинаются немецкое (Tschebyschev) и французское (Tchebychev) написания фамилии П.Л. Чебышёва, открывшего эти полиномы.

этой формулы на всей вещественной оси, а не только для значений аргумента $x \in [-1, 1]$.

Представление (2.27), в действительности, справедливо для любых $x \in \mathbb{R}$, если под $\operatorname{arccos} x$ понимать комплексное значение и, соответственно, рассматривать косинус от комплексного аргумента. Можно показать, что

$$T_n(x) = \operatorname{ch}(n \operatorname{arch} x), \quad (2.28)$$

где $\operatorname{ch} z = \frac{1}{2}(e^z + e^{-z})$ — гиперболический косинус, а arch — обратная к нему функция. Определение (2.28) удобно применять для вещественных аргументов x , таких что $|x| \geq 1$.

Предложение 2.3.1 *Функция $T_n(x)$, задаваемая формулой (2.27), — полином степени n , и его старший коэффициент при $n \geq 1$ равен 2^{n-1} .*

Доказательство. Мы проведём его индукцией по номеру n полинома Чебышёва. При $n = 0$ имеем $T_0(x) = 1$, при $n = 1$ справедливо $T_1(x) = x$, так что база индукции установлена.

Далее, из известной тригонометрической формулы

$$\cos \alpha + \cos \beta = 2 \cos \left(\frac{\alpha + \beta}{2} \right) \cos \left(\frac{\alpha - \beta}{2} \right)$$

следует, что

$$\begin{aligned} \cos((n+1) \arccos x) + \cos((n-1) \arccos x) \\ = 2 \cos(n \arccos x) \cos(\arccos x) \\ = 2x \cos(n \arccos x). \end{aligned}$$

Следовательно, в силу (2.27)

$$T_{n+1}(x) = 2x T_n(x) - T_{n-1}(x) \quad (2.29)$$

для любых $n = 1, 2, \dots$

Таким образом, если $T_{n-1}(x)$ и $T_n(x)$ являются полиномами степени $(n-1)$ и n соответственно, то $T_{n+1}(x)$ — также полином, степень которого на единицу выше степени $T_n(x)$, а старший коэффициент — в 2 раза больше. ■

Полученные в доказательстве рекуррентные формулы (2.29) позволяют последовательно выписывать явные алгебраические выражения

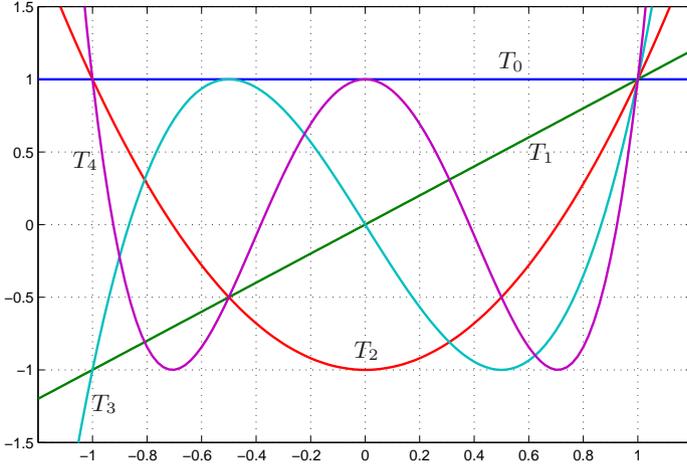


Рис. 2.6. Графики первых полиномов Чебышёва на интервале $[-1.2, 1.2]$.

для полиномов Чебышёва:

$$\begin{aligned}
 T_0(x) &= 1, \\
 T_1(x) &= x, \\
 T_2(x) &= 2x^2 - 1, \\
 T_3(x) &= 4x^3 - 3x, \\
 T_4(x) &= 8x^4 - 8x^2 + 1, \\
 T_5(x) &= 16x^5 - 20x^3 + 5x, \\
 &\dots
 \end{aligned}
 \tag{2.30}$$

По рекуррентным формулам (2.29) и следующим из них явным выражениям (2.30) полиномы Чебышёва единообразно определяются для любых значений аргумента x .

Рассмотрим кратко основные свойства полиномов Чебышёва. При чётном (нечётном) n полином Чебышёва $T_n(x)$ есть чётная (нечётная) функция от x . Действительно, выражение для $T_n(x)$ при чётном n содержит только чётные степени x (ноль считаем чётным числом), а при нечётном n — только нечётные степени x , что по индукции следует из рекуррентной формулы (2.29).

Найдём корни полиномов Чебышёва на вещественном интервале $[-1, 1]$. Исходя из представления (2.27), вспомним, каковы нули косинуса. Должно быть

$$n \arccos x = \frac{\pi}{2} + k\pi, \quad k \in \mathbb{Z},$$

причём в этой формуле k можно брать таким, чтобы область значений правой части не выходила за интервал $[0, n\pi]$, когда имеет смысл левая часть. Итак, корни полинома Чебышёва суть

$$\hat{x}_k = \cos \frac{(2k+1)\pi}{2n}, \quad k = 0, 1, \dots, n-1. \quad (2.31)$$

В целом, полином $T_n(x)$ имеет n вещественных различных корней, все они в самом деле находятся на интервале $[-1, 1]$ и выражаются в виде (2.31). Расположение корней полинома Чебышёва можно наглядно проиллюстрировать чертежом на Рис. 2.7, где эти корни соответствуют абсциссам точек пересечения единичной окружности с центром в начале координат с радиусами, откладываемыми через одинаковые доли развёрнутого угла. Из этой иллюстрации хорошо видно, что корни полинома Чебышёва расположены существенно неравномерно: они сгущаются к концам интервала $[-1, 1]$, а в середине более разрежены.

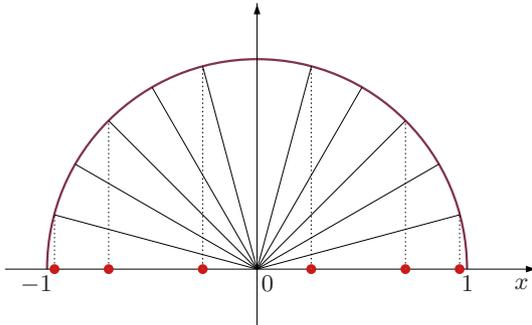


Рис. 2.7. Иллюстрация расположения корней полинома Чебышёва шестой степени.

Максимум модуля значений полинома Чебышёва на $[-1, 1]$ равен 1,

$$\max_{x \in [-1, 1]} |T_n(x)| = 1,$$

этот максимум достигается в точках $x_s = \cos(s\pi/n)$, $s = 0, 1, \dots, n$, причём $T_n(x_s) = (-1)^s$, $s = 0, 1, \dots, n$. Это непосредственно следует из тригонометрического определения полиномов Чебышёва (2.27), где внешний \cos должен достигать максимальных по модулю значений ± 1 в точках x_s , удовлетворяющих условию $n \arccos x = s\pi$, $s = 0, 1, \dots, n$.

Следующее свойство полиномов Чебышёва настолько важно, что мы оформим его как отдельное

Предложение 2.3.2 Среди полиномов степени n , $n \geq 1$, со старшим коэффициентом, равным 1, полином $\tilde{T}_n(x) := 2^{1-n} T_n(x)$ имеет на интервале $[-1, 1]$ наименьшее равномерное отклонение от нуля. Иными словами, если $Q_n(x)$ — полином степени n со старшим коэффициентом 1, то

$$\max_{x \in [-1, 1]} |Q_n(x)| \geq \max_{x \in [-1, 1]} |\tilde{T}_n(x)| = 2^{1-n}. \quad (2.32)$$

Доказательство. Предположим противное доказываемому, т. е. что для какого-то полинома $Q_n(x)$, имеющего старший коэффициент 1, выполняется неравенство

$$\max_{x \in [-1, 1]} |Q_n(x)| < \max_{x \in [-1, 1]} |\tilde{T}_n(x)|, \quad (2.33)$$

которое противоположно по смыслу неравенству (2.32). Тогда разность $(\tilde{T}_n(x) - Q_n(x))$ есть полином степени не выше $n - 1$. В то же время, в точках $x_s = \cos(s\pi/n)$, $s = 0, 1, \dots, n$, доставляющих полиному Чебышёва максимумы модуля на $[-1, 1]$, должно быть справедливо

$$\begin{aligned} \operatorname{sgn}(\tilde{T}_n(x_s) - Q_n(x_s)) &= \operatorname{sgn}((-1)^s 2^{1-n} - Q_n(x_s)) \\ &= \operatorname{sgn}((-1)^s 2^{1-n}) \quad \text{в силу (2.33)} \\ &= (-1)^s. \end{aligned}$$

Как следствие, на интервале $[x_s, x_{s+1}]$ полином $(\tilde{T}_n(x) - Q_n(x))$ меняет знак, и потому обязан иметь корень. Коль скоро это происходит для $s = 0, 1, \dots, n - 1$, т. е. всего n раз, то полином $(\tilde{T}_n(x) - Q_n(x))$ необходимо имеет n корней на $[-1, 1]$. Так как степень этого полинома не превосходит $n - 1$, полученные выводы можно примирить лишь при условии $(\tilde{T}_n(x) - Q_n(x)) = 0$, т. е. когда $Q_n(x) = \tilde{T}_n(x)$. Мы пришли к противоречию с допущением (2.33). \blacksquare

Полиномы $\tilde{T}_n(x)$, $n \geq 1$, фигурирующие в Предложении 2.3.2 и имеющие единичный старший коэффициент, называют *приведёнными полиномами Чебышёва*.

2.36 Применения полиномов Чебышёва

Доказательство Предложения 2.3.2 лишь косвенным образом использует то обстоятельство, что полиномы рассматриваются на интервале $[-1, 1]$. Фактически, мы опирались на свойство 3 полиномов Чебышёва достигать своих знакопеременных экстремумов в $n + 1$ точках этого интервала. Если в качестве области определения полиномов необходимо взять интервал $[a, b]$, отличный от $[-1, 1]$, то линейной заменой переменной

$$y = \frac{1}{2}(b + a) + \frac{1}{2}(b - a)x \quad (2.34)$$

интервал $[-1, 1]$ может быть преобразован в $[a, b]$. При этом обратное отображение $[a, b] \rightarrow [-1, 1]$ задаётся формулой

$$x = \frac{2y - (b + a)}{(b - a)}, \quad (2.35)$$

а корням полинома Чебышёва на $[-1, 1]$ соответствуют тогда точки

$$\hat{y}_k = \frac{1}{2}(b + a) + \frac{1}{2}(b - a) \cos \frac{(2k + 1)\pi}{2n}, \quad k = 0, 1, \dots, n - 1. \quad (2.36)$$

из интервала $[a, b]$. Свойство, аналогичное Предложению 2.3.2, будет верно на интервале $[a, b]$ для полинома, полученного из $T_n(x)$ с помощью линейной замены переменных (2.35).

Предложение 2.3.3 *Если $T_n(x)$ — n -ый полином Чебышёва, то полином переменной y , задаваемый как*

$$2^{1-2n} (b - a)^n \cdot T_n \left(\frac{2y - (b + a)}{b - a} \right) \quad (2.37)$$

имеет старший коэффициент 1 и на интервале $[a, b]$ равномерно наименее уклоняется от нуля среди всех полиномов степени n со старшим коэффициентом 1.

Доказательство. Первое утверждение Предложения следует из того, что при подстановке (2.35) в полиноме n -ой степени старший коэффициент приобретает дополнительный множитель $2^n / (b - a)^n$.

Из свойств полиномов Чебышёва следует, что на $[a, b]$ полином (2.37) достигает максимумов своего модуля, равных $2^{1-2n}(b-a)^n$, в точках

$$y_s = \frac{1}{2}(a+b) + \frac{1}{2}(b-a) \cos\left(\frac{s\pi}{n}\right),$$

$s = 0, 1, \dots, n$. Они получаются с помощью линейного преобразования (2.34) из аргументов $x_s = \cos(s\pi/n)$, $s = 0, 1, \dots, n$, доставляющих максимумы модуля полиному Чебышёва на $[-1, 1]$. Дальнейшие рассуждения повторяют доказательство Предложения 2.3.2, так как специфика интервала $[-1, 1]$ там, фактически, никак не использовалась. ■

Обратимся к поставленной в конце §2.2д задаче наиболее выгодного расположения узлов алгебраического интерполянта заданной степени n на некотором интервале $[a, b]$. Возьмём эти узлы корнями полинома (2.37), который получается в результате замены переменных (2.35) из полинома Чебышёва $(n+1)$ -ой степени $T_{n+1}(x)$, т. е. как

$$x_k = \frac{1}{2}(b+a) + \frac{1}{2}(b-a) \cos\left(\frac{(2k+1)\pi}{2(n+1)}\right), \quad k = 0, 1, \dots, n. \quad (2.38)$$

Тогда соответствующий полином

$$\omega_n(x) = (x-x_0)(x-x_1)\dots(x-x_n),$$

который фигурирует в формуле (2.24) для остаточного члена интерполяции, совпадёт с полиномом $(n+1)$ -ой степени вида (2.37). При этом $\omega_n(x)$ будет иметь наименьшее отклонение от нуля на $[a, b]$ в равномерной (чебышёвской) метрике (2.1), и в смысле этой метрики погрешность интерполирования при прочих равных условиях будет наименьшей возможной. Узлы интерполяции (2.38) называют *чебышёвскими узлами* на интервале $[a, b]$, а в совокупности они образуют *чебышёвскую сетку* на $[a, b]$.

Помимо интерполирования полиномы Чебышёва и их обобщения имеют и другие важные применения в различных задачах вычислительной математики и анализа [45, 57]. Очень важное значение имеют, к примеру, разложения функций в ряды по полиномам Чебышёва.

2.4 Алгебраическая интерполяция с кратными узлами

Кратным узлом называют, по определению, узел, в котором информация о функции задаётся более одного раза. Помимо значения функции это может быть какая-либо дополнительная информация о ней, например, значения производных и т. п. К задаче интерполяции с кратными узлами мы приходим, в частности, если степень интерполяционного полинома, который нужно однозначно построить по некоторым узлам, равна либо больше количества этих узлов.

Далее *задачей алгебраической интерполяции с кратными узлами* мы будем называть следующую постановку. Даны несовпадающие точки x_i , $i = 0, 1, \dots, n$, — узлы интерполирования, в которых заданы значения $y_i^{(k)}$, $k = 0, 1, \dots, N_i - 1$, — их принимают интерполируемая функция f и её производные $f^{(k)}(x)$. При этом число N_i называют *кратностью узла* x_i . Требуется построить полином $H_m(x)$ степени m , такой что

$$H_m^{(k)}(x_i) = y_i^{(k)}, \quad i = 0, 1, \dots, n, \quad k = 0, 1, \dots, N_i - 1. \quad (2.39)$$

Иными словами, в узлах x_i , $i = 0, 1, \dots, n$, как сам полином $H_m(x)$, так и все его производные $H_m^{(k)}(x)$ вплоть до заданных порядков N_i должны принимать предписанные им значения $y_i^{(k)}$.

Теорема 2.4.1 *Решение задачи алгебраической интерполяции с кратными узлами при $m = N_0 + N_1 + \dots + N_n - 1$ существует и единственно.*

Доказательство. В канонической форме полином $H_m(x)$ имеет вид

$$H_m(x) = \sum_{l=0}^m a_l x^l,$$

и для определения его коэффициентов a_0, a_1, \dots, a_m станем подставлять в него и в его производные $H_m'(x), H_m''(x), \dots$, аргументы x_i и использовать условия (2.39). Получим систему линейных алгебраических уравнений относительно a_0, a_1, \dots, a_m , в которой число уравнений $N_0 + N_1 + \dots + N_n$. При $m = N_0 + N_1 + \dots + N_n - 1$ оно совпадает с числом неизвестных, равным $m + 1$.

Обозначим получившуюся систему линейных уравнений как

$$Ga = y, \quad (2.40)$$

где G — квадратная $(m+1) \times (m+1)$ -матрица,

$a = (a_0, a_1, \dots, a_m)^\top \in \mathbb{R}^{m+1}$ — вектор неизвестных коэффициентов интерполяционного полинома,

$y = (y_0^{(0)}, y_0^{(1)}, \dots, y_0^{(N_0-1)}, y_1^{(0)}, y_1^{(1)}, \dots, y_n^{(N_n-1)})^\top \in \mathbb{R}^{m+1}$ — вектор, составленный из интерполяционных данных (2.39).

Матрица G зависит только от узлов x_0, x_1, \dots, x_n , и никак не зависит от данных $y_i^{(k)}$, $i = 0, 1, \dots, n$, $k = 0, 1, \dots, N_i - 1$. Хотя эту матрицу даже можно выписать в явном виде, её прямое исследование весьма сложно, и мы пойдём окольным путём.

Для определения свойств матрицы G рассмотрим однородную систему уравнений, отвечающую нулевой правой части $y = 0$, т. е.

$$Ga = 0.$$

Вектор правой части y образован значениями интерполируемой функции и её производных $y_i^{(k)}$ в узлах x_i , $i = 0, 1, \dots, n$. Однородная система $Ga = 0$ отвечает случаю $y_i^{(k)} = 0$ для всех $i = 0, 1, \dots, n$ и $k = 0, 1, \dots, N_i - 1$. Каким является вектор решений a этой системы?

Если правая часть в (2.40) — нулевая, то это означает, что полином $H_m(x)$ с учётом кратности имеет $N_0 + N_1 + \dots + N_n = m+1$ корней, т. е. больше, чем его степень. Следовательно, он necessarily равен нулю, а соответствующая однородная линейная система $Ga = 0$ имеет поэтому лишь нулевое решение $a = (a_0, a_1, \dots, a_n)^\top = (0, 0, \dots, 0)^\top$.

Итак, линейная комбинация столбцов матрицы G , равная нулю, может быть только тривиальной, т. е. с нулевыми коэффициентами. Следовательно, матрица G должна быть неособенной, а потому неоднородная линейная система (2.40) однозначно разрешима при любой правой части y . Это и требовалось доказать. ■

Задачу алгебраической интерполяции с кратными узлами в исследуемой нами постановке часто называют также задачей *эрмитовой интерполяции*, а сам полином $H_m(x)$ решающий эту задачу, называют *интерполяционным полиномом Эрмита*. Используемые при доказательстве Теоремы 2.4.1 рассуждения, в которых построение интерполя-

ционного полинома сводится к решению системы линейных алгебраических уравнений, носят конструктивный характер и позволяют практически решать задачу интерполяции с кратными узлами. Тем не менее, аналогично случаю интерполяции с простыми узлами, желательно иметь аналитическое решение в виде обозримого конечного выражения для интерполянта. Он может иметь форму Лагранжа либо форму Ньютона (см. подробности, к примеру, [3, 56]). Укажем способ построения его лагранжевой формы.

Аналогично §2.2б, при фиксированном наборе узлов x_0, x_1, \dots, x_n результат решения рассматриваемой задачи интерполяции линейно зависит от значений $y_0^{(0)}, y_0^{(1)}, \dots, y_0^{(N_0-1)}, y_1^{(0)}, y_1^{(1)}, \dots, y_n^{(N_n-1)}$. Более точно, если полином $P(x)$ решает задачу интерполяции по значениям $y = (y_0^{(0)}, y_0^{(1)}, \dots, y_n^{(N_n-1)})$, а полином $Q(x)$ решает задачу интерполяции с теми же узлами по значениям $z = (z_0^{(0)}, z_0^{(1)}, \dots, z_n^{(N_n-1)})$, то для любых вещественных чисел α и β полином $\alpha P(x) + \beta Q(x)$ решает задачу интерполяции для значений $\alpha y + \beta z = (\alpha y_0^{(0)} + \beta z_0^{(0)}, \alpha y_0^{(1)} + \beta z_0^{(1)}, \dots, \alpha y_n^{(N_n-1)} + \beta z_n^{(N_n-1)})$ на той же совокупности узлов.

Отмеченное свойство можно также усмотреть из выписанного при доказательстве Теоремы 2.4.1 представления вектора коэффициентов $a = (a_0, a_1, \dots, a_n)^T$ интерполяционного полинома как решения системы (2.40). Из него следует, что $a = G^{-1}y$, т.е. a линейно зависит от вектора данных y , образованного значениями $y_i^{(k)}$, $k = 0, 1, \dots, N_i - 1$, $i = 0, 1, \dots, n$.

Итак, свойством линейности можно воспользоваться для решения задачи интерполяции с кратными узлами «по частям», которые удовлетворяют отдельным интерполяционным условиям, а затем собрать эти части воедино. Иными словами, как и в случае интерполирования с простыми узлами, можно представить $H_m(x)$ в виде линейной комбинации

$$H_m(x) = \sum_{i=0}^n \sum_{k=0}^{N_i-1} y_i^{(k)} \cdot \phi_{ik}(x),$$

где внешняя сумма берётся по узлам, внутренняя — по порядкам производной, а $\phi_{ik}(x)$ — специальные «базисные» полиномы степени m , удовлетворяющие условиям

$$\phi_{ik}^{(l)}(x_j) = \begin{cases} 0, & \text{при } i \neq j \text{ или } l \neq k, \\ 1, & \text{при } i = j \text{ и } l = k. \end{cases} \quad (2.41)$$

У полинома $\phi_{ik}(x)$ в узле x_i не равна нулю лишь одна из производных, порядок которой k , тогда как производные всех других порядков зануляются в x_i . Кроме того, во всех остальных узлах полином $\phi_{ik}(x)$ и все его производные равны нулю. Фактически, полином $\phi_{ik}(x)$ отвечает линейной системе (2.40) с вектор-столбцом правой части y вида $(0, \dots, 0, 1, 0, \dots, 0)^\top$, в котором все элементы нулевые за исключением одного.

Каков конкретный вид этих базисных полиномов $\phi_{ik}(x)$? Перепишем условия (2.41) в виде

$$\phi_{ik}^{(l)}(x_i) = \delta_{kl}, \quad k = 0, 1, \dots, N_i - 1, \quad (2.42)$$

$$\begin{aligned} \phi_{ik}^{(l)}(x_j) &= 0, & j &= 0, 1, \dots, i-1, i+1, \dots, n, \\ & & l &= 0, 1, \dots, N_i - 1. \end{aligned} \quad (2.43)$$

Из второго условия следует, что

$$\phi_{ik}(x) = (x - x_0)^{N_0} \dots (x - x_{i-1})^{N_{i-1}} (x - x_{i+1})^{N_{i+1}} \dots (x - x_n)^{N_n} Q_{ik}(x),$$

где $Q_{ik}(x)$ — полином степени $N_i - 1$. Для его определения привлечём первое условие (2.42). И так далее.

Мы не будем завершать этого построения, так как дальнейшие выкладки весьма громоздки, а алгоритм нахождения полинома из приведённых рассуждения вполне ясен ...

Какова погрешность алгебраической интерполяции с кратными узлами?

Теорема 2.4.2 Пусть $f \in C^{m+1}[a, b]$, т. е. функция f непрерывно дифференцируема $m + 1$ раз на интервале $[a, b]$. Погрешность $R_m(f, x)$ её интерполирования по попарно различным узлам $x_0, x_1, \dots, x_n \in [a, b]$ с кратностями N_0, N_1, \dots, N_n полиномом $H_m(x)$ степени m при условии $m = N_0 + N_1 + \dots + N_n - 1$ может быть представлена в виде

$$R_m(f, x) = f(x) - H_m(x) = \frac{f^{(m+1)}(\xi(x))}{(m+1)!} \cdot \prod_{i=0}^n (x - x_i)^{N_i}, \quad (2.44)$$

где $\xi(x)$ — некоторая точка из $]a, b[$, зависящая от x .

Доказательство. Обозначим для удобства через $\Omega(x)$ произведение линейных множителей со степенями из правой части равенства (2.44),

т. е.

$$\Omega(x) := \prod_{i=0}^n (x - x_i)^{N_i}.$$

Это — аналог функции $\omega_n(x)$, введённой в §2.2д и широко используемой выше.

Если $x = x_i$ для одного из узлов интерполирования, $i = 0, 1, \dots, n$, то $R_m(f, x) = 0$, но в то же время и $\Omega(x) = 0$. Поэтому в (2.44) в качестве ξ можно взять любую точку из открытого интервала $]a, b[$.

Предположим теперь, что точка x из интервала интерполирования $[a, b]$ не совпадает ни с одним из узлов x_i , $i = 0, 1, \dots, n$. Введём вспомогательную функцию

$$\psi(z) := f(z) - H_m(z) - K\Omega(z),$$

где K — числовая константа, равная

$$K = \frac{f(x) - H_m(x)}{\Omega(x)}.$$

Функция $\psi(z)$ имеет нули в узлах x_0, x_1, \dots, x_n и, кроме того, по построению обращается в нуль в точке x , так что общее число нулей этой функции равно $n + 2$. На основании теоремы Ролля можно заключить, что производная $\psi'(z)$ обращается в нуль по крайней мере в $n + 1$ точках, расположенных в интервалах между x, x_1, \dots, x_n . Но в узлах x_0, x_1, \dots, x_n функция $\psi(z)$ имеет нули с кратностями N_0, N_1, \dots, N_n соответственно, и потому в них производная $\psi'(z)$ имеет нули кратности $N_0 - 1, N_1 - 1, \dots, N_n - 1$ (нулевая кратность означает отсутствие нуля в узле). Таким образом, всего эта производная $\psi'(z)$ имеет с учётом кратности $(N_0 + N_1 + \dots + N_n - n - 1) + (n + 1) = m + 1$ нулей на $[a, b]$.

Продолжая аналогичные рассуждения, получим, что вторая производная $\psi''(z)$ будет иметь с учётом кратности по крайней мере m нулей на интервале $[a, b]$ и т. д. При каждом последующем дифференцировании нули у производных функции $\psi(z)$ могут возникать или исчезать, но их суммарная кратность уменьшается всякий раз на единицу. Наконец, $(m + 1)$ -ая производная зануляется на $[a, b]$ хотя бы один раз.

Итак, на интервале $[a, b]$ обязательно найдётся по крайней мере одна точка ξ , такая что $\psi^{(m+1)}(\xi) = 0$. Но

$$\psi^{(m+1)}(z) = f^{(m+1)}(z) - K(m+1)!,$$

поскольку $H_m(x)$ — полином степени m и $H_m^{(m+1)}(z)$ зануляется, а $\Omega(z)$ есть многочлен степени $m+1$ со старшим коэффициентом 1. Итак,

$$K = \frac{f^{(m+1)}(\xi)}{(m+1)!}$$

для некоторой точки ξ , зависящей от x . Этим завершается доказательство теоремы. \blacksquare

Отметим, что при наличии одного узла кратности m интерполяционный полином Эрмита должен совпасть с полиномом Тейлора, а формула (2.44) превращается в известную формулу остаточного члена для полинома Тейлора. Если же все узлы интерполяции простые, то (2.44) совпадает с полученной ранее формулой погрешности простой интерполяции (2.24).

2.5 Общие факты алгебраической интерполяции

Как с теоретической, так и с практической точек зрения интересен вопрос о том, насколько малой может быть сделана погрешность интерполирования при возрастании числа узлов. Вообще, имеет ли место сходимость интерполяционных полиномов к интерполируемой функции при неограниченном возрастании количества узлов?

Чтобы строго сформулировать соответствующие вопросы и общие результаты о сходимости алгебраических интерполянтов, необходимо формализовать некоторые понятия.

Определение 2.5.1 Пусть для интервала $[a, b]$ задана бесконечная треугольная матрица узлов

$$\begin{pmatrix} x_0^{(1)} & 0 & 0 & 0 & \cdots \\ x_0^{(2)} & x_1^{(2)} & 0 & 0 & \cdots \\ x_0^{(3)} & x_1^{(3)} & x_2^{(3)} & 0 & \cdots \\ \vdots & \vdots & \vdots & \ddots & \ddots \end{pmatrix},$$

такая что в каждой её строке расположены различные точки интервала $[a, b]$, т. е. $x_i^{(n)} \in [a, b]$ для всех положительных целых n и любых

$i = 0, 1, \dots, n$, причём $x_i^{(n)} \neq x_j^{(n)}$ для $i \neq j$. Говорят, что на интервале $[a, b]$ задан интерполяционный процесс, если элементы n -ой строки этой матрицы берутся в качестве узлов интерполяции, по которым строится последовательность интерполянтов $g_n(x)$, $n = 1, 2, \dots$.

Если все интерполянты $g_n(x)$ являются алгебраическими полиномами, то будем употреблять термин *алгебраический интерполяционный процесс*.

Определение 2.5.2 *Интерполяционный процесс для функции f называется сходящимся в точке $y \in [a, b]$, если последовательность значений интерполянтов $g_n(y) \rightarrow f(y)$ при $n \rightarrow \infty$. Интерполяционный процесс для функции f на интервале $[a, b]$, порождающий последовательность интерполянтов $g_n(x)$, называется сходящимся равномерно, если $\max_{x \in [a, b]} |f(x) - g_n(x)| \rightarrow 0$ при $n \rightarrow \infty$.*

Отметим, что помимо равномерной сходимости интерполяционного процесса, когда отклонение одной функции от другой измеряется в равномерной (чебышёвской) метрике (2.1), иногда необходимо рассматривать сходимость в других смыслах. Например, это может быть среднеквадратичная сходимость, задаваемая метрикой (2.3), или ещё какая-нибудь другая.

Определённую уверенность в положительном ответе на поставленные в начале параграфа вопросы даёт известная из математического анализа

Теорема Вейерштрасса о равномерном приближении.

Если $f : [a, b] \rightarrow \mathbb{R}$ — непрерывная функция, то для всякого $\epsilon > 0$ существует полином $P_n(x)$ степени $n = n(\epsilon)$, равномерно приближающий функцию f с погрешностью, не превосходящей ϵ , т. е. такой, что

$$\max_{x \in [a, b]} |f(x) - P_n(x)| \leq \epsilon.$$

Этот результат служит теоретической основой равномерного приближения непрерывных функций алгебраическими полиномами, обеспечивая существование полинома, который сколь угодно близок к заданной непрерывной функции в смысле расстояния (2.1). Вместе с тем, теорема Вейерштрасса слишком обща и не даёт ответа на конкретные

вопросы о решении задачи интерполирования, где требуется совпадение значений функции и её интерполянта на данном множестве точек-узлов.

Как следует из результатов §2.2д и §2.3 огромное влияние на погрешность интерполяции оказывает расположение узлов. В частности, рассмотренные в §2.3 чебышёвские сетки являются наилучшими возможными в условиях, когда неизвестна какая-либо дополнительная информация об интерполируемой функции.

Что касается равномерных сеток, то для них один из первых примеров расходимости интерполяционных процессов привёл в 1910 году С.Н. Бернштейн, рассмотрев на интервале $[-1, 1]$ алгебраическую интерполяцию функции $f(x) = |x|$ по равноотстоящим узлам, включая и концы этого интервала. Не слишком трудными рассуждениями показывается (см. [7, 25]), что с возрастанием числа узлов соответствующий интерполяционный полином не стремится к $|x|$ ни в одной точке интервала $[-1, 1]$, отличной от $-1, 0$ и 1 . Может показаться, что причиной плохой сходимости интерполяционного процесса в примере С.Н. Бернштейна является отсутствие гладкости интерполируемой функции, но это верно лишь отчасти.

Предположим, что интерполируемая функция f имеет бесконечную гладкость, т.е. $f \in C^\infty[a, b]$, и при этом её производные растут «не слишком быстро». В последнее условие будем вкладывать следующий смысл:

$$\sup_{x \in [a, b]} |f^{(n)}(x)| < M^n, \quad n = 0, 1, 2, \dots, \quad (2.45)$$

где M не зависит от n . Тогда из Теоремы 2.2.2 следует, что погрешность алгебраического интерполирования по n узлам может быть оценена сверху как

$$\frac{(M(b-a))^{n+1}}{(n+1)!},$$

то есть при $n \rightarrow \infty$ очевидно сходится к нулю вне зависимости от расположения узлов интерполяции. Иными словами, любой алгебраический интерполяционный процесс на интервале $[a, b]$ будет равномерно сходиться к такой функции f .

Условие (2.45) влечёт сходимость ряда Тейлора для функции f в любой точке из $[a, b]$, и такие функции называются *аналитическими* на рассматриваемом множестве [40]. Это очень важный, хотя и не слишком широкий класс функций, которые являются ближайшим обобщением алгебраических полиномов.

Но в самом общем случае при алгебраическом интерполировании бесконечно гладких функций погрешность всё-таки может не сходиться к нулю, даже при «вполне разумном» расположении узлов. По-видимому, наиболее известный пример такого рода привёл немецкий математик К. Рунге. В примере Рунге функция

$$\Upsilon(x) = \frac{1}{1 + 25x^2}$$

интерполируется алгебраическими полиномами на интервале $[-1, 1]$ с равномерным расположением узлов интерполяции $x_i = -1 + 2i/n$, $i = 0, 1, 2, \dots, n$. Оказывается, что тогда

$$\lim_{n \rightarrow \infty} \max_{x \in [-1, 1]} |\Upsilon(x) - P_n(x)| = \infty,$$

где $P_n(x)$ — интерполяционный полином n -ой степени. При этом с ростом n вблизи концов интервала интерполирования $[-1, 1]$ у полиномов $P_n(x)$ возникают сильные колебания (часто называемые также *осцилляциями*), размах которых стремится к бесконечности (см. Рис. 2.8). Получается, что хотя в узлах интерполирования значения функции $\Upsilon(x)$ совпадают со значениями интерполяционного полинома, между этим узлами $P_n(x)$ и $\Upsilon(x)$ могут отличаться сколь угодно сильно, даже несмотря на плавный (бесконечно гладкий) характер изменения функции $\Upsilon(x)$.

Интересно, что на интервале $[-\kappa, \kappa]$, где $\kappa \approx 0.726$, рассматриваемый интерполяционный процесс равномерно сходится к $\Upsilon(x)$ (см. [64]). Кроме того, полезно отметить, что функция $\Upsilon(x)$ имеет производные всех порядков для любого вещественного аргумента x , но у концов интервала интерполирования $[-1, 1]$ эти производные растут очень быстро и уже не удовлетворяют условию (2.45). Таким образом, несмотря на простой вид, функция $\Upsilon(x)$ из примера Рунге своим поведением слишком непохожа на полиномы, производные от которых растут умеренно и, начиная с некоторого порядка, исчезают. Эти интересные вопросы относятся уже к теории функций.

Что касается чебышёвских сеток, то они обеспечивают равномерную сходимость интерполяционного процесса к функциям, которые подчиняются так называемому условию Дини-Липшица [8, 25, 50]. Это очень слабое условие, которому заведомо удовлетворяют большинство встречающихся на практике функций. Практичным достаточным условием, при котором выполняется условие Дини-Липшица, является *обоб-*

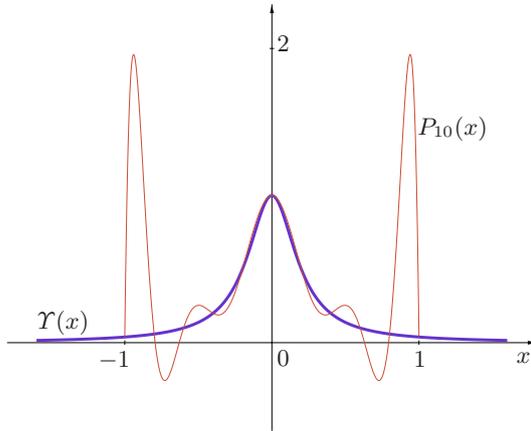


Рис. 2.8. Интерполяция полиномом 10-й степени в примере Рунге

ущённое условие Липшица: для любых x, y из области определения функции

$$|f(x) - f(y)| \leq C |x - y|^\alpha, \quad (2.46)$$

для некоторых C и $0 < \alpha \leq 1$. Иными словами, справедлива

Теорема 2.5.1 *Если узлами интерполирования берутся чебышёвские сетки, то интерполяционный процесс сходится равномерно для любой функции, удовлетворяющей обобщённому условию Липшица (2.46).*

Обоснование этого утверждения можно увидеть, к примеру, в [25]. Тем не менее, для общих непрерывных функций имеет место следующий отрицательный результат:

Теорема Фабера [7, 8, 56] *Не существует бесконечной треугольной матрицы узлов из заданного интервала, такой что соответствующий ей алгебраический интерполяционный процесс сходился бы равномерно для любой непрерывной функции на этом интервале.*

В частности, даже для интерполяционного процесса по узлам полиномов Чебышёва существуют примеры непрерывных функций, для которых этот алгебраический интерполяционный процесс всюду расходится. Подробности можно найти в [25].

Но отрицательный результат теоремы Фабера характеризует, скорее, слишком большую общность математического понятия непрерывной функции, которая может оказаться слишком необычной и «неудобной» и не похожей на то, что мы интуитивно вкладываем в смысл «непрерывности». Здесь уместно напомнить примеры непрерывных нигде не дифференцируемых функций (примеры Вейерштрасса или ван дер Ваардена). Столь же экзотичен пример непрерывной функции, к которой не сходится равномерно интерполяционный процесс по чебышёвским сеткам. Что касается теоремы Фабера, то она утверждает лишь то, что класс непрерывных в классическом смысле функций является слишком широким, чтобы для него существовал один (или даже несколько) интерполяционных процессов, обеспечивающих равномерную сходимость для любой функции.

Чересчур большая общность понятия непрерывной функции была осознана математиками почти сразу после своего появления, в первой половине XIX века. Она стимулировала работы по формулировке дополнительных естественных условий, которые выделяли бы классы функций, непрерывных в более сильных смыслах, которые позволяли бы свободно выполнять те или иные традиционные операции (например, взятие производной почти всюду в области определения и т. п.). Именно эти причины вызвали появление условий Липшица, Дини-Липшица, обобщённого условия Липшица и ряда других им аналогичных.

Следует отметить, что для общих непрерывных функций имеет место «оптимистичный» результат, имеющий, правда, небольшую практическую ценность:

Теорема Марцинкевича [7, 8, 56] *Для любой непрерывной на заданном интервале функции f найдётся такая бесконечная треугольная матрица узлов из этого интервала, что соответствующий ей алгебраический интерполяционный процесс для функции f сходится равномерно.*

Интересно отметить, что ситуация со сходимостью интерполяционных процессов в среднеквадратичном смысле более благоприятна. Согласно результату, который получили П. Эрдёш и П. Туран (см. [25]), для любой положительной весовой функции существует треугольная матрица узлов из интервала интерполирования, по которой интерполяционный процесс будет сходиться. Иными словами, равномерная схо-

димось предъявляет к функции требования, более сильные чем среднеквадратичная сходимость.

Ещё один вывод из представленных выше примеров и результатов заключается в том, что алгебраические полиномы, несмотря на определённые удобства работы с ними, оказываются довольно капризным инструментом интерполирования достаточно общих непрерывных и даже гладких функций. Как следствие, нам нужно иметь более гибкие инструменты интерполяции. Их развитию и будут посвящены следующие параграфы.

2.6 Сплайны

2.6а Элементы теории

Пусть заданный интервал $[a, b]$ разбит на подинтервалы $[x_{i-1}, x_i]$, $i = 1, 2, \dots, n$, так что $a = x_0$ и $x_n = b$. *Сплайном* на $[a, b]$ называется функция, которая вместе со своими производными вплоть до некоторого порядка непрерывна на всём интервале $[a, b]$, и на каждом подинтервале $[x_{i-1}, x_i]$ является полиномом. Максимальная на всём интервале $[a, b]$ степень полиномов, задающих сплайн, называется *степенью сплайна*. Разность между степенью сплайна и наивысшим порядком его производной, которая непрерывна на $[a, b]$, называется *дефектом сплайна*. Наконец, точки x_i , $i = 0, 1, \dots, n$, — концы подинтервалов $[x_{i-1}, x_i]$ — называют *узлами сплайна*.

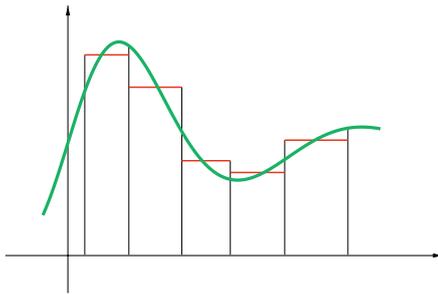


Рис. 2.9. Кусочно-постоянная интерполяция функции.

Термин «сплайн» является удачным заимствованием из английско-

го языка, где слово «spline» означает гибкую (обычно стальную) линейку, которую, изгибая, использовали чертёжники для проведения гладкой линии между данными фиксированными точками. В середине XX века этот термин вошёл в математику и перекочевал во многие языки мира.

Почему именно кусочные полиномы? К идее их введения можно прийти, к примеру, с помощью следующих неформальных мотиваций. Содержательный (механический, физический, биологический и т. п.) смысл имеют, как правило, производные порядка не выше 2–4, и именно их мы можем видеть в различных математических моделях реальных явлений.⁶ Пятые производные — это уже экзотика, а производные шестого и более высоких порядков при описании реальности не встречаются. В частности, разрывы производных высоких «нефизических» порядков у функции никак не ощутимы. Поэтому для сложно изменяющихся производных высоких порядков необходимые «нужные» значения в фиксированных узлах можно назначить, к примеру, с помощью простейшей кусочно-постоянной или кусочно-линейной интерполяции, а затем добиться желаемой гладкости исходной функции с помощью последовательного применения нескольких операций интегрирования. Так получается кусочно-полиномиальная функция.

Понятие «сплайн-функции» было введено И. Шёнбергом в 1946 году [73], хотя различные применения тех объектов, которые впоследствии были названы «сплайнами», встречались в математике на протяжении предшествующей сотни лет. Пионером здесь следует назвать, по-видимому, Н.И. Лобачевского, который в статье [70] явно использовал конструкции сплайнов и так называемых *B*-сплайнов.⁷

С середины XX века по настоящее время сплайны нашли широкие применения в математике и её приложениях. В вычислительных технологиях они могут использоваться, в частности, для приближения функций, при решении дифференциальных и интегральных уравнений. Если сплайн применяется для решения задачи интерполяции, то он называется *интерполяционным*. Другими словами, интерполяционный сплайн — это сплайн, принимающий в заданных точках \tilde{x}_i , $i = 0, 1,$

⁶Характерный пример: в книге А.К. Маловичко, О.Л. Тарунина «Использование высших производных при обработке и интерпретации результатов геофизических наблюдений» (Москва, издательство «Недра», 1981 год) рассматриваются производные второго и третьего порядков.

⁷Вклад Н.И. Лобачевского даже дал повод некоторым авторам назвать специальный вид сплайнов *сплайнами Лобачевского*.

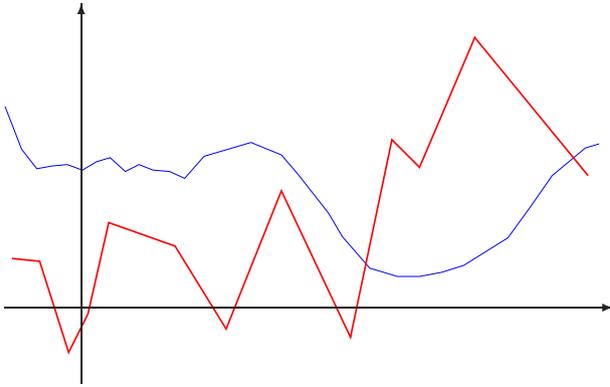


Рис. 2.10. Простейшие сплайны — кусочно-линейные функции.

\dots, r , — узлах интерполяции — требуемые значения y_i . Эти узлы интерполяции, вообще говоря, могут не совпадать с узлами сплайна x_i , $i = 0, 1, \dots, n$, задающим интервалы полиномиальности.

Так как степень полинома равна наивысшему порядку его ненулевой производной, то сплайны дефекта нуль — это функции задаваемые на всём интервале $[a, b]$ одной полиномиальной формулой. Таким образом, термин «дефект» весьма точно выражает то, сколько сплайну «не хватает» до полноценного полинома. С другой стороны, именно наличие дефекта обеспечивает сплайну большую гибкость в сравнении с полиномами и делает сплайны в некоторых ситуациях более удобным инструментом приближения и интерполирования функций. Чем больше дефект сплайна, тем больше он отличается от полинома и тем более специфичны его свойства. Но слишком большой дефект приводит к существенному понижению общей гладкости сплайна. В значительном числе приложений сплайнов вполне достаточным оказываются сплайны с минимально возможным дефектом 1, и только такие сплайны мы будем рассматривать далее в нашей книге.

Простейший «настоящий» сплайн имеет дефект 1 и степень 1, будучи «непрерывно склеенным» в узлах x_i , $i = 1, 2, \dots, n-1$. Иными словами, это — кусочно-линейная функция, имеющая, несмотря на свою простоту, богатые приложения в математике.⁸ Сплайны второй степени

⁸Вспомним, к примеру, «ломанные Эйлера», которые применяются при доказательстве существования решения задачи Коши для обыкновенных дифференци-

часто называют *параболическими сплайнами*.

Далее мы будем рассматривать интерполяционные сплайны, узлы которых x_0, x_1, \dots, x_n совпадают с узлами интерполирования.

Если степень сплайна равна d , то для его полного определения необходимо знать $n(d+1)$ значений коэффициентов полиномов, задающих сплайн на n подинтервалах $[x_{i-1}, x_i]$, $i = 1, 2, \dots, n$. В то же время, в случае дефекта 1 имеется

$d(n-1)$ условий непрерывности самого сплайна и его производных вплоть до $(d-1)$ -го порядка в узлах x_1, x_2, \dots, x_{n-1} ,

$(n+1)$ условие интерполяции в узлах x_0, x_1, \dots, x_n .

Всего $d(n-1) + (n+1) = n(d+1) - (d-1)$ штук, и потому для определения сплайна не хватает $d-1$ условий, которые обычно задают дополнительно на концах интервала $[a, b]$.

Сказанное имеет следующие важные следствия. Если решать задачу интерполяции с помощью сплайна чётной степени, требуя, чтобы на каждом подинтервале $[x_{i-1}, x_i]$ сплайн являлся бы полиномом чётной степени, то число $(d-1)$ подлежащих доопределению параметров оказывается нечётным. Поэтому на одном из концов интервала $[a, b]$ приходится налагать больше условий, чем на другом. Это приводит, во-первых, к асимметрии задачи, и, во-вторых, может вызвать неустойчивость при определении параметров сплайна. Наконец, интерполяционный сплайн чётной степени при некоторых естественных краевых условиях (периодических, к примеру) может просто не существовать.

Отмеченные недостатки могут решаться, в частности, выбором узлов сплайна отличными от узлов интерполяции. Мы далее не будем останавливаться на преодолении этих затруднений и рассмотрим интерполяционные сплайны нечётной степени 3.

2.66 Интерполяционные кубические сплайны

Наиболее популярны в вычислительной математике сплайны третьей степени с дефектом 1, называемые также *кубическими сплайнами*. Эту популярность можно объяснить относительной простотой этих сплайнов и тем обстоятельством, что они вполне достаточны для от-

альных уравнений [37].

слеживания непрерывности вторых производных функций, что необходимо, например, во многих законах механики и физики.

Пусть задан набор узлов $x_0, x_1, \dots, x_n \in [a, b]$, который, как и прежде, мы называем *сеткой*. Величину $h_i = x_i - x_{i-1}$, $i = 1, 2, \dots, n$, назовём *шагом сетки*. Кубический интерполяционный сплайн на интервале $[a, b]$ с сеткой $a = x_0 < x_1 < \dots < x_n = b$, узлы которой являются также узлами интерполяции — это функция $S(x)$, удовлетворяющая следующим условиям:

- 1) $S(x)$ — полином третьей степени на каждом из подинтервалов $[x_{i-1}, x_i]$, $i = 1, 2, \dots, n$;
- 2) $S(x) \in C^2[a, b]$;
- 3) $S(x_i) = y_i$, $i = 0, 1, 2, \dots, n$.

Для построения такого сплайна $S(x)$ нужно определить $4n$ неизвестных величин — по 4 коэффициента полинома третьей степени на каждом из n штук подинтервалов $[x_{i-1}, x_i]$, $i = 1, 2, \dots, n$.

В нашем распоряжении имеются

$3(n - 1)$ условий непрерывности самой функции $S(x)$, её первой и второй производных во внутренних узлах x_1, x_2, \dots, x_{n-1} ;

$(n + 1)$ условие интерполяции $S(x_i) = y_i$, $i = 0, 1, 2, \dots, n$.

Таким образом, для определения $4n$ неизвестных величин мы имеем всего $(4n - 2)$ условий. Два недостающих условия определяют различными способами, среди которых часто используются, к примеру, такие:

$$(I) \quad S'(a) = \beta_0, \quad S'(b) = \beta_n,$$

$$(II) \quad S''(a) = \gamma_0, \quad S''(b) = \gamma_n,$$

$$(III) \quad S^{(k)}(a) = S^{(k)}(b), \quad k = 0, 1, 2,$$

где $\beta_0, \beta_n, \gamma_0, \gamma_n$ — данные вещественные числа. Условия (I) и (II) соответствуют заданию на концах интервала $[a, b]$ первой или второй производной искомого сплайна, а условие (III) — это условие периодического продолжения сплайна с интервала $[a, b]$ на более широкое подмножество вещественной оси.

Мы рассмотрим подробно случай (II) задания краевых условий:

$$S''(a) = S''(x_0) = \gamma_0,$$

$$S''(b) = S''(x_n) = \gamma_n.$$

Будем искать кусочно-полиномиальное представление нашего кубического сплайна в специальном виде, привязанном к узлам сплайна x_i : пусть

$$S(x) = \alpha_{i-1} + \beta_{i-1}(x - x_{i-1}) + \gamma_{i-1} \frac{(x - x_{i-1})^2}{2} + \vartheta_{i-1} \frac{(x - x_{i-1})^3}{6} \quad (2.47)$$

для $x \in [x_{i-1}, x_i]$, $i = 1, 2, \dots, n$. Ясно, что в такой форме представления сплайна величины β_0 и γ_0 совпадают по смыслу с теми, что даются в условиях (I)–(II) выше. Более того, из представления (2.47) вытекает, что

$$S''(x_{i-1}) = \gamma_{i-1}, \quad i = 1, 2, \dots, n.$$

Вторая производная $S''(x)$ является линейной функцией на $[x_{i-1}, x_i]$, и с учётом (2.47) должно быть:

$$S''(x) = \gamma_{i-1} + \vartheta_{i-1}(x - x_{i-1}), \quad x \in [x_{i-1}, x_i]. \quad (2.48)$$

С другой стороны, вид этой линейной функции полностью определяется двумя её крайними значениями γ_{i-1} и γ_i на концах подинтервала $[x_{i-1}, x_i]$. Поэтому вместо (2.48) можно выписать более определённое представление:

$$S''(x) = \gamma_{i-1} \frac{x_i - x}{h_i} + \gamma_i \frac{x - x_{i-1}}{h_i}$$

для $x \in [x_{i-1}, x_i]$, $i = 1, 2, \dots, n$. В этих формулах при $i = 0$ и $i = n$ мы задействуем известные нам из условия (II) значения γ_0 и γ_n второй производной S'' на левом и правом концах интервала $[a, b]$. Очевидно, что построенная таким образом функция $S''(x)$ удовлетворяет условию «непрерывной склейки» в узлах x_1, x_2, \dots, x_{n-1} , т. е.

$$S''(x_i - 0) = S''(x_i + 0), \quad i = 1, 2, \dots, n - 1.$$

Чтобы восстановить S по S'' , нужно теперь взять дважды первообразную (неопределённый интеграл) от $S''(x)$. Проинтегрировав, получим

$$S(x) = \gamma_{i-1} \frac{(x_i - x)^3}{6h_i} + \gamma_i \frac{(x - x_{i-1})^3}{6h_i} + C_1 x + C_2 \quad (2.49)$$

с какими-то константами C_1 и C_2 . Но нам будет удобно представить это выражение в несколько другом виде:

$$S(x) = \gamma_{i-1} \frac{(x_i - x)^3}{6h_i} + \gamma_i \frac{(x - x_{i-1})^3}{6h_i} + K_1(x_i - x) + K_2(x - x_{i-1}), \quad (2.50)$$

где K_1 и K_2 — также константы.⁹ Насколько законен переход к такой форме? Из сравнения (2.49) и (2.50) следует, что C_1 и C_2 должны быть связаны с K_1 и K_2 посредством формул

$$\begin{aligned} C_1 &= -K_1 + K_2, \\ C_2 &= K_1 x_i - K_2 x_{i-1}. \end{aligned}$$

У выписанной системы линейных уравнений относительно K_1 и K_2 определитель равен $x_{i-1} - x_i = -h_i$, и он не зануляется. Поэтому переход от C_1 и C_2 к K_1 и K_2 — это неособенная замена переменных. Следовательно, оба представления (2.49) и (2.50) совершенно равносильны друг другу.

Для определения K_1 и K_2 воспользуемся интерполяционными условиями. Подставляя в выражение (2.50) значения $x = x_{i-1}$ и используя условия $S(x_{i-1}) = y_{i-1}$, $i = 1, 2, \dots, n$, будем иметь

$$\gamma_{i-1} \frac{(x_i - x_{i-1})^3}{6h_i} + K_1(x_i - x_{i-1}) = y_{i-1},$$

т. е.

$$\gamma_{i-1} \frac{h_i^2}{6} + K_1 h_i = y_{i-1},$$

откуда

$$K_1 = \frac{y_{i-1}}{h_i} - \frac{\gamma_{i-1} h_i}{6}.$$

Совершенно аналогичным образом, подставляя в (2.50) значение $x = x_i$ и используя условие $S(x_i) = y_i$, найдём

$$K_2 = \frac{y_i}{h_i} - \frac{\gamma_i h_i}{6}.$$

⁹Строго говоря, константы C_1 , C_2 , K_1 , K_2 нужно было бы снабдить ещё дополнительным индексом i , показывающими их зависимость от подинтервала $[x_{i-1}, x_i]$, к которому они относятся. Мы не делаем этого ради краткости изложения.

Выражение сплайна на подинтервале $[x_{i-1}, x_i]$, $i = 1, 2, \dots, n$, выглядит поэтому следующим образом:

$$S(x) = y_{i-1} \frac{x_i - x}{h_i} + y_i \frac{x - x_{i-1}}{h_i} + \gamma_{i-1} \frac{(x_i - x)^3 - h_i^2(x_i - x)}{6h_i} + \gamma_i \frac{(x - x_{i-1})^3 - h_i^2(x - x_{i-1})}{6h_i}. \quad (2.51)$$

Оно не содержит уже величин α_i , β_i и ϑ_i , которые фигурировали в исходном представлении (2.47) для $S(x)$, но неизвестными остались γ_1 , $\gamma_2, \dots, \gamma_{n-1}$.

Чтобы завершить определение вида сплайна, т. е. найти $\gamma_1, \gamma_2, \dots, \gamma_{n-1}$, можно воспользоваться условием непрерывности первой производной $S'(x)$ в узлах x_1, x_2, \dots, x_{n-1} :

$$S'(x_i - 0) = S'(x_i + 0), \quad i = 1, 2, \dots, n - 1. \quad (2.52)$$

Продифференцировав по x формулу (2.51), получим для $x \in [x_{i-1}, x_i]$

$$S'(x) = \frac{y_i - y_{i-1}}{h_i} - \gamma_{i-1} \frac{3(x_i - x)^2 - h_i^2}{6h_i} + \gamma_i \frac{3(x - x_{i-1})^2 - h_i^2}{6h_i}. \quad (2.53)$$

Следовательно, с учётом того, что $x_i - x_{i-1} = h_i$,

$$\begin{aligned} S'(x_i) &= \frac{y_i - y_{i-1}}{h_i} + \gamma_{i-1} \frac{h_i^2}{6h_i} + \gamma_i \frac{3(x_i - x_{i-1})^2 - h_i^2}{6h_i} \\ &= \frac{y_i - y_{i-1}}{h_i} + \gamma_{i-1} \frac{h_i}{6} + \gamma_i \frac{h_i}{3}. \end{aligned} \quad (2.54)$$

С другой стороны, сдвигая все индексы в (2.53) на единицу вперёд, получим для подинтервала $x \in [x_i, x_{i+1}]$ представление

$$S'(x) = \frac{y_{i+1} - y_i}{h_{i+1}} - \gamma_i \frac{3(x_{i+1} - x)^2 - h_{i+1}^2}{6h_{i+1}} + \gamma_{i+1} \frac{3(x - x_i)^2 - h_{i+1}^2}{6h_{i+1}}.$$

Следовательно, с учётом того, что $x_{i+1} - x_i = h_{i+1}$,

$$\begin{aligned} S'(x_i) &= \frac{y_{i+1} - y_i}{h_{i+1}} - \gamma_i \frac{3(x_{i+1} - x_i)^2 - h_{i+1}^2}{6h_{i+1}} - \gamma_{i+1} \frac{h_{i+1}^2}{6h_{i+1}} \\ &= \frac{y_{i+1} - y_i}{h_i} - \gamma_i \frac{h_{i+1}}{3} - \gamma_{i+1} \frac{h_{i+1}}{6}. \end{aligned} \quad (2.55)$$

Приравнивание, согласно (2.52), производных (2.54) и (2.55), которые получены в узлах x_i с соседних подинтервалов $[x_{i-1}, x_i]$ и $[x_i, x_{i+1}]$, приводит к соотношениям

$$\begin{cases} \frac{h_i}{6} \gamma_{i-1} + \frac{h_i + h_{i+1}}{3} \gamma_i + \frac{h_{i+1}}{6} \gamma_{i+1} = \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i}, \\ i = 1, 2, \dots, n-1, \\ \gamma_0 \text{ и } \gamma_n \text{ заданы.} \end{cases} \quad (2.56)$$

Это система линейных алгебраических уравнений относительно неизвестных переменных $\gamma_1, \gamma_2, \dots, \gamma_{n-1}$, имеющая матрицу

$$\frac{1}{6} \begin{pmatrix} 2(h_1 + h_2) & h_2 & & & 0 \\ h_2 & 2(h_2 + h_3) & h_3 & & & \\ & h_3 & 2(h_3 + h_4) & \ddots & & \\ & & \ddots & \ddots & \ddots & \\ 0 & & & & h_{n-1} & 2(h_{n-1} + h_n) \end{pmatrix},$$

в которой ненулевыми являются лишь три диагонали — главная и соседние с ней — сверху и снизу. Такие матрицы называются *трёхдиагональными*. (см. §3.8). Кроме того, эта матрица обладает свойством *диагонального преобладания* (стр. 223): стоящие на её главной диагонали элементы $\frac{1}{3}(h_i + h_{i+1})$ по модулю больше, чем сумма модулей внедиагональных элементов этой же строки. В силу признака Адамара (он рассматривается и обосновывается в §3.2e) такие матрицы неособенны. Как следствие, система линейных уравнений (2.56) относительно γ_i , $i = 1, 2, \dots, n-1$, однозначно разрешима при любой правой части, а искомый сплайн всегда существует и единствен. Для нахождения решения системы (2.56) с трёхдиагональной матрицей может быть с успехом применён метод прогонки, описываемый ниже в §3.8.

Интересен вопрос о погрешности интерполирования функций и их производных с помощью кубических сплайнов, и ответ на него даёт следующая

Теорема 2.6.1 Пусть $f(x) \in C^p[a, b]$, $p = 1, 2, 3, 4$, а $S(x)$ — интерполяционный кубический сплайн с краевыми условиями (I), (II) или (III), построенный по значениям $f(x)$ на сетке $a = x_0 < x_1 < \dots < x_n = b$ из интервала $[a, b]$, с шагом $h_i = x_i - x_{i-1}$, $i = 1, 2, \dots, n$, причём узлы

интерполяции являются также узлами сплайна. Тогда для $k = 0, 1, 2$ и $k \leq p$ справедливо соотношение

$$\max_{x \in [a, b]} |f^{(k)}(x) - S^{(k)}(x)| = O(h^{p-k}),$$

где $h = \max_i h_i$.

При формулировке этого утверждения и далее в этой книге мы пользуемся символом $O(\cdot)$ — «О-большое», введённым Э. Ландау и широко используемым в современной математике и её приложениях. Для двух переменных величин u и v пишут, что $u = O(v)$, если отношение u/v есть величина ограниченная в рассматриваемом процессе. В формулировке Теоремы 2.6.1 и в других ситуациях, где идёт речь о шаге сетки h , мы всюду имеем в виду $h \rightarrow 0$. Удобство использования символа $O(\cdot)$ состоит в том, что, показывая качественный характер зависимости, он не требует явного выписывания констант, которые должны фигурировать в соответствующих отношениях.

Обоснование Теоремы 2.6.1 разбивается на ряд частных случаев, соответствующих различным значениям гладкости p и порядка производной k . Их доказательства можно увидеть, к примеру, в [10, 12, 29]. Повышение гладкости p интерполируемой функции $f(x)$ выше, чем $p = 4$, уже не оказывает влияния на погрешность интерполирования, так как интерполяционный сплайн кубический, т. е. имеет степень 3. С другой стороны, свои особенности имеет также случай $p = 0$, когда интерполируемая функция всего лишь непрерывна, и мы не приводим здесь полную формулировку соответствующего результата о погрешности.

Отметим, что, в отличие от алгебраических интерполянтов, последовательность интерполяционных кубических сплайнов на равномерной сетке узлов всегда сходится к интерполируемой непрерывной функции. Это относится, в частности, и к функции $\mathcal{T}(x) = 1/(1 + 25x^2)$ из примера Рунге (см. §2.5). Важно также, что с повышением гладкости интерполируемой функции до определённого предела сходимость эта улучшается.

С другой стороны, интерполирование сплайнами иллюстрирует также интересное явление *насыщения* численных методов, когда, начиная с какого-то порядка, увеличение гладкости исходных данных задачи уже не приводит к увеличению точности результата. Соответствующие численные методы называют *насыщаемыми*. Напротив, *ненасыщаемые численные методы*, там, где их удаётся построить и применить, дают всё более точное решение при увеличении гладкости решения [37].

Основной недостаток понятий насыщаемости / ненасыщаемости состоит в трудности практического определения гладкости данных, которые присутствуют в предъявленной к решению задаче.

2.6в Экстремальное свойство кубических сплайнов

Интерполяционные кубические сплайны $S(x)$, удовлетворяющие на концах рассматриваемого интервала $[a, b]$ дополнительным условиям

$$S''(a) = S''(b) = 0, \quad (2.57)$$

называются *естественными* или *натуральными сплайнами*. Их замечательное свойство состоит в том, что они минимизируют функционал

$$\mathcal{E}(f) = \int_a^b (f''(x))^2 dx,$$

выражающий в первом приближении энергию упругой деформации гибкой стальной линейки, форма которой описывается функцией $f(x)$ на интервале $[a, b]$. Краевые условия (2.57) соответствуют при этом линейке, свободно закреплённой на концах.

Как известно, потенциальная энергия изгибания малого участка упругого тела пропорциональна квадрату его кривизны (скорости изгибания в зависимости от длины дуги) в данной точке. Кривизна плоской кривой, задаваемой уравнением $y = f(x)$, равна, как известно,

$$\frac{f''(x)}{(1 + (f'(x))^2)^{3/2}}$$

(см., к примеру, [34, 58]). Поэтому упругая энергия однородной линейки, принимающей форму кривой $y = f(x)$ на интервале $[a, b]$, при условии приблизительного постоянства $f'(x)$, пропорциональна

$$\int_a^b (f''(x))^2 dx.$$

Теорема 2.6.2 *Если $S(x)$ — естественный сплайн, построенный по узлам $a = x_0 < x_1 < \dots < x_n = b$, а $\varphi(x)$ — любая другая дважды гладкая функция, принимающая в этих узлах те же значения, что и $S(x)$, то $\mathcal{E}(\varphi) \geq \mathcal{E}(S)$, причём неравенство строго для $\varphi \neq S$.*

Доказательство этого факта не очень сложно и может быть найдено, к примеру, в [2, 10, 32].

Будучи предоставленной самой себе, упругая линейка, закреплённая в узлах интерполирования, принимает форму, которая, как известно из физики, должна минимизировать энергию своей упругой деформации. Таким образом, эта форма очень близка к кубическому сплайну.

Сформулированное свойство называют *экстремальным свойством* естественных сплайнов,¹⁰ и оно служит началом большого и важного направления в теории сплайнов.

2.7 Нелинейные методы интерполяции

Рассмотренные нами выше методы интерполяции (в частности, алгебраической), были *линейными* в том смысле, что результат решения задачи интерполяции при фиксированных узлах линейно зависел от данных. При этом класс интерполирующих функций \mathcal{S} образует линейное векторное пространство над полем \mathbb{R} : любая линейная комбинация функций также является функцией заданного вида, решающей задачу интерполяции для линейной комбинации данных. Но существуют и другие, нелинейные, методы интерполирования, для которых не выполнено сформулированное выше свойство. Эти методы также широко применяются при практической интерполяции, так как обладают многими важными достоинствами.

Нелинейными называют методы интерполяции, в которых класс интерполирующих функций \mathcal{S} не является линейным векторным пространством. Важнейший частный случай нелинейных методов интерполяции — это интерполяция с помощью рациональных функций вида

$$y = y(x) = \frac{a_0 + a_1x + a_2x^2 + \dots}{b_0 + b_1x + b_2x^2 + \dots}. \quad (2.58)$$

Итак, пусть в узлах x_0, x_1, \dots, x_n заданы значения функции y_0, y_1, \dots, y_n . Нам нужно найти рациональную дробь вида (2.58), такую что $y_i = y(x_i), i = 0, 1, \dots, n$.

Поскольку дробь не меняется от умножения числителя и знаменателя на одно и то же ненулевое число, то для какого-нибудь одного из коэффициентов a_i или b_i может быть выбрано произвольное наперёд заданное значение. Кроме того, параметры a_i и b_i должны удовлетворять

¹⁰Иногда также говорят о *вариационном свойстве* естественных сплайнов.

$n + 1$ условиям интерполяции в узлах, так что всего этих параметров мы можем извлечь из условия задачи $(n + 2)$ штук. Этим ограничением определяется общее число неизвестных параметров, т.е. сумма степеней многочленов числителя и знаменателя в дроби (2.58).

Представление (2.58) равносильно тождеству

$$a_0 - b_0y + a_1x - b_1xy + a_2x^2 - b_2x^2y + \dots = 0. \quad (2.59)$$

Коль скоро при $x = x_i$ должно быть $y = y_i$, $i = 0, 1, \dots, n$, то получаем ещё $(n + 1)$ числовых равенств

$$a_0 - b_0y_i + a_1x_i - b_1x_iy_i + a_2x_i^2 - b_2x_i^2y_i + \dots = 0, \quad (2.60)$$

$i = 0, 1, \dots, n$. Соотношения (2.59)–(2.60) можно трактовать, как условия линейной зависимости с коэффициентами $a_0, -b_0, a_1, -b_1, \dots$ для вектор-столбцов

$$\begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \begin{pmatrix} y \\ y_0 \\ y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \begin{pmatrix} x \\ x_0 \\ x_1 \\ \vdots \\ x_n \end{pmatrix}, \quad \begin{pmatrix} xy \\ x_0y_0 \\ x_1y_1 \\ \vdots \\ x_ny_n \end{pmatrix}, \quad \begin{pmatrix} x^2 \\ x_0^2 \\ x_1^2 \\ \vdots \\ x_n^2 \end{pmatrix}, \quad \begin{pmatrix} x^2y \\ x_0^2y_0 \\ x_1^2y_1 \\ \vdots \\ x_n^2y_n \end{pmatrix}, \quad \dots$$

размера $(n + 2)$. Как следствие, определитель

$$\det \begin{pmatrix} 1 & y & x & xy & x^2 & x^2y & \dots \\ 1 & y_0 & x_0 & x_0y_0 & x_0^2 & x_0^2y_0 & \dots \\ 1 & y_1 & x_1 & x_1y_1 & x_1^2 & x_1^2y_1 & \dots \\ 1 & y_2 & x_2 & x_2y_2 & x_2^2 & x_2^2y_2 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \\ 1 & y_n & x_n & x_ny_n & x_n^2 & x_n^2y_n & \dots \end{pmatrix}$$

составленной из этих столбцов матрицы размера $(n + 2) \times (n + 2)$ должен быть равен нулю. Разлагая его по первой строке и разрешая полученное равенство нулю относительно y , мы действительно получим выражение для y в виде отношения двух многочленов от x .

Реализация описанного выше приёма требует нахождения значений определителя числовых $(n + 1) \times (n + 1)$ -матриц, и далее в §3.13 мы

рассмотрим соответствующие методы. Отметим, что в популярных системах компьютерной математики Scilab, MATLAB, Maple, Mathematica и др. для этого существует готовая встроенная функция `det`.

Пример 2.7.1 Рассмотрим в качестве примера интерполяцию дробно-рациональной функцией таблицы значений



2.8 Численное дифференцирование

Дифференцированием называется, как известно, процесс нахождения производной от заданной функции или же численного значения этой производной в заданной точке. Необходимость выполнения дифференцирования возникает весьма часто и вызвана огромным распространением этой операции в современной математике и её приложениях. Производная бывает нужна и сама по себе, как мгновенная скорость тех или иных процессов, и как вспомогательное средство для построения более сложных процедур, например, в методе Ньютона для численного решения уравнений и систем уравнений (см. §§4.4г и 4.5б).

В настоящее время наиболее распространены три следующих способа вычисления производных:

- символьное (аналитическое) дифференцирование,
- численное дифференцирование,
- алгоритмическое (автоматическое) дифференцирование.

Символьным (аналитическим) дифференцированием называют процесс построения по функции, задаваемой каким-то выражением, производной функции, основываясь на известных из математического анализа правилах дифференцирования составных функций (суммы, разности, произведения, частного, композиции, обратной функции и т. п.) и известных производных для простейших функций. Основы символьного (аналитического) дифференцирования являются предметом математического анализа (точнее, дифференциального исчисления), а более продвинутые результаты по этой теме входят в курсы компьютерной алгебры.

На принципах, похожих на символьное (аналитическое) дифференцирование, основывается алгоритмическое (автоматическое) дифференцирование, но при этом оперируют не выражениями для производных,

а их численными значениями при данных значениях аргументов функции. Как символьное (аналитическое) дифференцирование, так и алгоритмическое (автоматическое) дифференцирование требуют знания выражения для функции или хотя бы компьютерной программы для её вычисления. Мы кратко рассмотрим алгоритмическое дифференцирование в §2.9.

Численным дифференцированием называется процесс нахождения значения производной от функции, использующий значения этой функции в некотором наборе точек её области определения. Таким образом, если функция задана таблично, т.е. лишь на конечном множестве значений аргумента, либо процедура определения значений этой функции не может быть выписана в виде выражения или детерминированной программы, то альтернатив численному дифференцированию нет. В частности, иногда в виде такого «чёрного ящика» мы вынуждены представлять вычисление значений функции, аналитическое выражение для которой существует, но является слишком сложным или неудобным для дифференцирования первыми двумя способами.

В основе методов численного дифференцирования лежат различные идеи. Самая первая состоит в том, чтобы доопределить (восстановить) таблично заданную функцию до функции непрерывного аргумента, к которой уже применима обычная операция дифференцирования. Теория интерполирования, которой посвящены предшествующие параграфы, может оказаться в высшей степени полезной при реализации такого подхода. В частности, таблично заданную функцию можно заменить её интерполяционным полиномом, и его производные считать производными рассматриваемой функции. Для этого годится также интерполяция сплайнами или какими-либо другими функциями, а в целом описанный выше подход к численному дифференцированию называют *интерполяционным подходом*.

2.8a Интерполяционный подход

Итак, пусть задан набор узлов $x_0, x_1, \dots, x_n \in [a, b]$, т.е. сетка с шагом $h_i = x_i - x_{i-1}$, $i = 1, 2, \dots, n$. Кроме того, заданы значения функции f_0, f_1, \dots, f_n , такие что $f_i = f(x_i)$, $i = 0, 1, \dots, n$. Ниже мы рассмотрим простейший вариант интерполяционного подхода, в котором используется алгебраическая интерполяция.

Начнём со случая, когда применяется интерполяционный полином первой степени, который мы строим по двум соседним узлам сетки, т.е.

по x_{i-1} и x_i , $i = 1, 2, \dots, n$:

$$\begin{aligned} P_{1,i}(x) &= \frac{x - x_i}{x_{i-1} - x_i} f_{i-1} + \frac{x - x_{i-1}}{x_i - x_{i-1}} f_i \\ &= \frac{f_i - f_{i-1}}{x_i - x_{i-1}} x + \frac{f_{i-1}x_i - f_i x_{i-1}}{x_i - x_{i-1}}, \end{aligned}$$

где у интерполяционного полинома добавлен дополнительный индекс « i », указывающий на ту пару узлов, по которым он построен. Поэтому производная равна

$$P'_{1,i}(x) = \frac{f_i - f_{i-1}}{x_i - x_{i-1}} = \frac{f_i - f_{i-1}}{h_i}.$$

Это значение можно взять за приближение к производной от рассматриваемой функции на интервале $]x_{i-1}, x_i[$, $i = 1, 2, \dots, n$.

Во внутренних узлах сетки — x_1, x_2, \dots, x_{n-1} , — т. е. там, где встречаются два подинтервала, производную можно брать по любой из возможных формул

$$f'(x_i) \approx f_{\bar{x},i} := \frac{f_i - f_{i-1}}{x_i - x_{i-1}} \quad \text{— разделённая разность назад,} \quad (2.61)$$

$$f'(x_i) \approx f_{x,i} := \frac{f_{i+1} - f_i}{x_{i+1} - x_i} \quad \text{— разделённая разность вперёд.} \quad (2.62)$$

Обе они примерно равнозначны и выбор конкретной из них может быть делом соглашения, удобства или целесообразности. Например, от направления этой разности может решающим образом зависеть устойчивость разностных схем для численного решения дифференциальных уравнений.

Построим теперь интерполяционные полиномы Лагранжа второй степени по трём соседним точкам сетки x_{i-1}, x_i, x_{i+1} , $i = 1, 2, \dots, n-1$.

Имеем

$$\begin{aligned}
 P_{2,i}(x) &= \frac{(x-x_i)(x-x_{i+1})}{(x_{i-1}-x_i)(x_{i-1}-x_{i+1})} f_{i-1} + \frac{(x-x_{i-1})(x-x_{i+1})}{(x_i-x_{i-1})(x_i-x_{i+1})} f_i \\
 &\quad + \frac{(x-x_{i-1})(x-x_i)}{(x_{i+1}-x_{i-1})(x_{i+1}-x_{i-1})} f_{i+1} \\
 &= \frac{x^2 - (x_i + x_{i+1})x + x_i x_{i+1}}{(x_{i-1} - x_i)(x_{i-1} - x_{i+1})} f_{i-1} \\
 &\quad + \frac{x^2 - (x_{i-1} + x_{i+1})x + x_{i-1} x_{i+1}}{(x_i - x_{i-1})(x_i - x_{i+1})} f_i \\
 &\quad + \frac{x^2 - (x_{i-1} + x_i)x + x_{i-1} x_i}{(x_{i+1} - x_{i-1})(x_{i+1} - x_i)} f_{i+1}.
 \end{aligned}$$

Поэтому

$$\begin{aligned}
 P'_{2,i}(x) &= \frac{2x - (x_i + x_{i+1})}{(x_{i-1} - x_i)(x_{i-1} - x_{i+1})} f_{i-1} + \frac{2x - (x_{i-1} + x_{i+1})}{(x_i - x_{i-1})(x_i - x_{i+1})} f_i \\
 &\quad + \frac{2x - (x_{i-1} + x_i)}{(x_{i+1} - x_{i-1})(x_{i+1} - x_i)} f_{i+1}.
 \end{aligned}$$

Воспользуемся теперь тем, что $x_i - x_{i-1} = h_i$, $x_{i+1} - x_i = h_{i+1}$. Тогда $x_{i+1} - x_{i-1} = h_i + h_{i+1}$, а результат предшествующих выкладок может быть записан в виде

$$\begin{aligned}
 f'(x) \approx P'_{2,i}(x) &= \frac{2x - x_i - x_{i+1}}{h_i(h_i + h_{i+1})} f_{i-1} \\
 &\quad - \frac{2x - x_{i-1} - x_{i+1}}{h_i h_{i+1}} f_i + \frac{2x - x_{i-1} - x_i}{h_{i+1}(h_i + h_{i+1})} f_{i+1}.
 \end{aligned} \tag{2.63}$$

Формула (2.63) может применяться при вычислении значения производной в произвольной точке x для случая общей неравномерной сетки. Предположим теперь для простоты, что сетка равномерна, т. е. $h_i = h = \text{const}$, $i = 1, 2, \dots, n$. Кроме того, для таблично заданной функции на практике обычно наиболее интересны производные в тех же точках, где задана сама функция, т. е. в узлах x_0, x_1, \dots, x_n . В точке $x = x_i$ из (2.63) получаем для первой производной формулу

$$f'(x_i) \approx f_{\bar{x},i} = \frac{f_{i+1} - f_{i-1}}{2h}, \tag{2.64}$$

называемую *формулой центральной разности*. Подставляя в (2.63) аргумент $x = x_{i-1}$ и сдвигая в получающемся результате индекс на +1, получим

$$f'(x_i) \approx \frac{-3f_i + 4f_{i+1} - f_{i+2}}{2h}.$$

Подставляя в (2.63) аргумент $x = x_{i+1}$ и сдвигая в получающемся результате индекс на (-1), получим

$$f'(x_i) \approx \frac{f_{i-2} - 4f_{i-1} + 3f_i}{2h}.$$

Займёмся теперь выводом формул для второй производной. Используя интерполяционный полином второй степени, можно найти:

$$f''(x_i) \approx P''_{2,i}(x) = \frac{2}{h_i(h_i + h_{i+1})} f_{i-1} - \frac{2}{h_i h_{i+1}} f_i + \frac{2}{h_{i+1}(h_i + h_{i+1})} f_{i+1}.$$

В частности, на равномерной сетке с $h_i = h = \text{const}$, $i = 1, 2, \dots, n$ имеем

$$f''(x_i) \approx \frac{f_{i-1} - 2f_i + f_{i+1}}{h^2}. \quad (2.65)$$

Эта формула широко используется в вычислительной математике, и по аналогии с (2.61)–(2.62) часто обозначается кратко как $f_{x\bar{x}}$. Естественно, что полученные выражения для второй производной не зависят от аргумента x .

Несмотря на то, что проведённые выше рассуждения основывались на применении интерполяционного полинома Лагранжа, для взятия производных произвольных порядков на сетке общего вида удобнее использовать интерполяционный полином Ньютона, в котором члены являются полиномами возрастающих степеней.

Выпишем ещё без вывода формулы численного дифференцирования на равномерной сетке, полученные по четырём точкам, т. е. с применением интерполяционного полинома третьей степени: для первой

производной —

$$f'(x_i) \approx \frac{1}{6h}(-11f_i + 18f_{i+1} - 9f_{i+2} + 2f_{i+3}), \quad (2.66)$$

$$f'(x_i) \approx \frac{1}{6h}(-2f_{i-1} - 3f_i + 6f_{i+1} - f_{i+2}), \quad (2.67)$$

$$f'(x_i) \approx \frac{1}{6h}(f_{i-2} - 6f_{i-1} + 3f_i + 2f_{i+1}), \quad (2.68)$$

$$f'(x_i) \approx \frac{1}{6h}(-2f_{i-3} + 9f_{i-2} - 18f_{i-1} + 11f_i), \quad (2.69)$$

для второй производной —

$$f''(x_i) \approx \frac{1}{h^2}(2f_i - 5f_{i+1} + 4f_{i+2} - f_{i+3}), \quad (2.70)$$

$$f''(x_i) \approx \frac{1}{h^2}(f_{i-1} - 2f_i + f_{i+1}), \quad (2.71)$$

$$f''(x_i) \approx \frac{1}{h^2}(-f_{i-3} + 4f_{i-2} - 5f_{i-1} + 2f_i). \quad (2.72)$$

В формуле (2.71) один из четырёх узлов, по которым строилась формула, никак не используется, а сама формула совпадает с формулой (2.65), полученной по трём точкам. Отметим красивую двойственность формул (2.66) и (2.69), (2.67) и (2.68), а также (2.70) и (2.72). Неслучаен также тот факт, что сумма коэффициентов при значениях функции в узлах во всех формулах равна нулю: он является следствием того, что производная постоянной функции — нуль.



Рис. 2.11. Шаблон формулы второй разностной производной (2.65).

В связи с численным дифференцированием и во многих других вопросах вычислительной математики чрезвычайно полезно понятие шаблона (сеточной) формулы, под которым мы будем понимать совокупность охватываемых этой формулой узлов сетки. Более точно, *шаблон формулы* численного дифференцирования — это множество узлов

сетки, входящих в правую часть этой формулы, явным образом либо в качестве аргументов используемых значений функции. Например, шаблоном формулы (2.65) для вычисления второй производной на равномерной сетке

$$f''(x_i) \approx \frac{f_{i-1} - 2f_i + f_{i+1}}{h^2}$$

являются три точки x_{i-1} , x_i , x_{i+1} (см. Рис. 2.11), в которых должны быть заданы f_{i-1} , f_i , f_{i+1} . Особенно разнообразны формы шаблонов в случае двух и более независимых переменных.

2.8б Оценка погрешности численного дифференцирования

Пусть для численного нахождения k -ой производной функции применяется формула численного дифференцирования Φ , имеющая шаблон Ψ и использующая значения функции в узлах этого шаблона. Если $f(x)$ — дифференцируемая необходимое число раз функция, такая что $f_i = f(x_i)$ для всех узлов $x_i \in \Psi$, то какова может быть погрешность вычисления $f^{(k)}(x)$ по формуле Φ ? Вопрос этот можно адресовать как к целому интервалу значений аргумента, так и локально, только к той точке x_i , которая служит аргументом левой части формулы численного дифференцирования.

Если рассматриваемая формула выведена в рамках интерполяционного подхода, то заманчивой идеей является получение ответа прямым дифференцированием полученных нами ранее выражений (2.23) и (2.24) для погрешности интерполирования. Этот путь оказывается очень непростым, так как применение, к примеру, выражения (2.24) требует достаточной гладкости функции $\xi(x)$, о которой мы можем сказать немного. Даже если эта гладкость имеется у $\xi(x)$, полученные оценки будут содержать производные $\xi'(x)$ и пр., о которых мы знаем ещё меньше. Наконец, шаблон некоторых формул численного дифференцирования содержит меньше точек, чем это необходимо для построения интерполяционных полиномов нужной степени. Такова, к примеру, формула «центральной разности» для первой производной или формула для второй производной (2.71), построенная по четырём точкам на основе полинома 3-й степени. Тем не менее, явные выражения для остаточного члена формул численного дифференцирования на этом пути можно получить методом, который напоминает вывод фор-

мулы для погрешности алгебраического интерполирования. Подробности изложены, к примеру, в книгах [17, 56].

Рассмотрим ниже детальнее более простой и достаточно универсальный способ оценивания погрешностей, основанный на разложениях по формуле Тейлора. Суть этого способа заключается, во-первых, в выписывании по формуле Тейлора разложений для функций, входящих в правую часть формулы численного дифференцирования, и, во-вторых, в аккуратном учёте членов этих разложений с целью получить, по возможности, наиболее точное выражение для ошибки.

Поясним эту методику на примере оценки погрешности для формулы «центральной разности» (2.64):

$$f'(x_i) \approx f_{\bar{x},i} = \frac{f_{i+1} - f_{i-1}}{2h}.$$

Предположим, что $f \in C^3[x_{i-1}, x_{i+1}]$, т. е. функция f трижды непрерывно дифференцируема на интервале между узлами формулы. Подставляя её в (2.64) и разлагая относительно точки x_i по формуле Тейлора с остаточным членом в форме Лагранжа вплоть до членов второго порядка, получим

$$\begin{aligned} f_{\bar{x},i} &= \frac{1}{2h} \left(\left(f(x_i) + hf'(x_i) + \frac{h^2}{2} f''(x_i) + \frac{h^3}{6} f'''(\xi_+) \right) \right. \\ &\quad \left. - \left(f(x_i) - hf'(x_i) + \frac{h^2}{2} f''(x_i) - \frac{h^3}{6} f'''(\xi_-) \right) \right) \\ &= f'(x_i) + \frac{h^2}{12} f'''(\xi_+) + \frac{h^2}{12} f'''(\xi_-), \end{aligned}$$

где ξ_+ и ξ_- — некоторые точки из открытого интервала $]x_{i-1}, x_{i+1}[$. Поэтому

$$f_{\bar{x},i} - f'(x_i) = \frac{h^2}{12} (f'''(\xi_+) + f'''(\xi_-)) = \frac{\alpha h^2}{6},$$

где $\alpha = \frac{1}{2}(f'''(\xi_+) + f'''(\xi_-))$. В целом справедлива оценка

$$|f_{\bar{x},i} - f'(x_i)| \leq \frac{M_3}{6} h^2,$$

в которой $M_3 = \max_{\xi} |f'''(\xi)|$ для $\xi \in]x_{i-1}, x_{i+1}[$. То есть, на трижды непрерывно дифференцируемых функциях погрешность вычисления производной по формуле «центральной разности» равна $O(h^2)$ для равномерной сетки шага h .

Определение 2.8.1 *Станем говорить, что приближённая формула (численного дифференцирования, интегрирования и т. п.) или же приближённый численный метод имеют p -ый порядок точности (или порядок аппроксимации), если на равномерной сетке с шагом h их погрешность является величиной $O(h^p)$, т. е. не превосходит Ch^p , где C — константа, не зависящая от h .*

Нередко понятие порядка точности распространяют и на неравномерные сетки, в которых шаг h_i меняется от узла к узлу. Тогда роль величины h играет какой-нибудь «характерный размер», описывающий данную сетку, например, $h = \max_i h_i$. Порядок точности — важная количественная мера погрешности формулы или метода, и при прочих равных условиях более предпочтительной является та формула или тот метод, которые имеют более высокий порядок точности. Но следует чётко осознавать, что порядок точности имеет асимптотический характер и отражает поведение погрешности при стремлении шагов сетки к нулю. Если этого стремления нет и шаг сетки остаётся «достаточно большим», то вполне возможны ситуации, когда метод меньшего порядка точности даёт лучшие результаты, поскольку множитель при h^p в оценке погрешности у него меньше.

Другое необходимое замечание состоит в том, что понятие порядка формулы или метода основывается на сравнении скорости убывания погрешности со скоростью убывания степенных функций $1, x, x^2, \dots, x^k, \dots$, то есть существенно завязано на степенную шкалу. Иногда (не слишком часто) эта шкала оказывается не вполне адекватной реальному поведению погрешности.

Пример 2.8.1 Пусть на вещественной оси задана равномерная сетка шага h , включающая в себя узлы $0, \pm h, \pm 2h$ и т. д. Для функции $y = g(x)$ рассмотрим интерполяцию значения $g(0)$ полусуммой

$$\frac{1}{2}(g(-h) + g(h)), \quad (2.73)$$

т. е. простейшим интерполяционным полиномом первой степени по узлам $-h$ и h . Каков будет порядок погрешности такой интерполяции в

зависимости от h для различных функций $g(x)$?¹¹

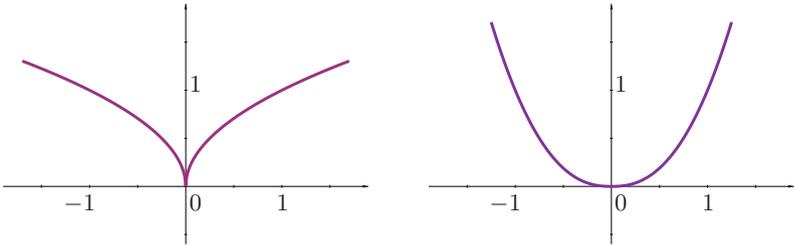


Рис. 2.12. Графики функции $y = |x|^\alpha$ при $0 < \alpha < 1$ и $\alpha > 1$.

Для функции $g(x) = |x|^\alpha$, $\alpha > 0$, погрешность интерполяции будет, очевидно, равна h^α , так что её порядок равен α . Он может быть нецелым числом (в частности, дробным) и даже сколь угодно малым.

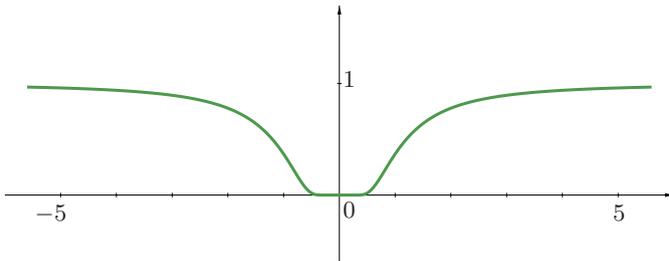


Рис. 2.13. График функции $y = \exp(-1/x^2)$.

Возьмём в качестве $g(x)$ функцию

$$g(x) = \begin{cases} \exp\left(-\frac{1}{x^2}\right), & \text{при } x \neq 0, \\ 0, & \text{при } x = 0, \end{cases}$$

известную в математическом анализе как пример бесконечно гладкой, но не аналитической (т. е. не разлагающейся в степенной ряд) функции. Погрешность интерполяции значения этой функции в нуле с помощью формулы (2.73) равна $\exp(-1/h^2)$, при $h \rightarrow 0$ она убывает быстрее любой степени h , так что порядок точности нашей интерполяции оказывается бесконечно большим. Но такой же бесконечно большой порядок точности интерполирования будет демонстрировать здесь функция

¹¹Идея этого примера заимствована из пособия [42], задача 4.2.

$y = x^2 g(x)$, хотя для неё погрешность $h^2 \exp(-1/h^2)$ убывает существенно быстрее. ■

Из выкладок, проведённых для определения погрешности формулы «центральной разности», хорошо видна особенность метода разложений по формуле Тейлора: его *локальный* характер, вытекающий из свойств самой формулы Тейлора. Наши построения оказываются «привязанными» к определённому узлу (или узлам) сетки, относительно которого и следует строить все разложения, чтобы обеспечить взаимные уничтожения их ненужных членов. Как следствие, в этом специальном узле (узлах) мы можем быстро оценить погрешность. Но за пределами этого узла (узлов), в частности, между узлами сетки всё гораздо сложнее и не так красиво, поскольку взаимные уничтожения членов могут уже не происходить.

Какой порядок точности имеют другие формулы численного дифференцирования?

Методом разложений по формуле Тейлора для дважды гладкой функции f нетрудно получить оценки

$$|f_{x,i} - f'(x_i)| \leq \frac{M_2}{2} h, \quad |f_{\bar{x},i} - f'(x_i)| \leq \frac{M_2}{2} h, \quad (2.74)$$

где $M_2 = \max_{\xi} |f''(\xi)|$ по ξ из соответствующего интервала между узлами. Таким образом, разность вперёд (2.61) и разность назад (2.61) имеют всего лишь первый порядок точности. Отметим, что для дважды непрерывно дифференцируемых функций оценки (2.74) уже не могут быть улучшены и достигаются, к примеру, на функции $f(x) = x^2$.

Конспективно изложим другие результаты о точности формул численного дифференцирования:

$$f'(x_i) = \frac{1}{2h} (-3f_i + 4f_{i+1} - f_{i+2}) + O(h^2),$$

$$f'(x_i) = \frac{1}{2h} (f_{i-2} - 4f_{i-1} + 3f_i) + O(h^2),$$

$$f'(x_i) = \frac{1}{6h} (-2f_{i-1} - 3f_i + 6f_{i+1} - f_{i+2}) + O(h^3),$$

$$f'(x_i) = \frac{1}{6h} (f_{i-2} - 6f_{i-1} + 3f_i + 2f_{i+1}) + O(h^3).$$

Оценим теперь погрешность формулы (2.65) для второй производной

$$f''(x_i) \approx f_{x\bar{x},i} = \frac{f_{i-1} - 2f_i + f_{i+1}}{h^2}.$$

Обозначая для краткости $f'_i = f'(x_i)$ и $f''_i = f''(x_i)$, получим

$$\begin{aligned} f_{x\bar{x},i} &= \frac{1}{h^2} \left(\left(f_i - hf'_i + \frac{h^2}{2} f''_i - \frac{h^3}{6} f'''_i + \frac{h^4}{24} f^{(4)}(\xi_-) \right) - 2f_i \right. \\ &\quad \left. + \left(f_i + hf'_i + \frac{h^2}{2} f''_i + \frac{h^3}{6} f'''_i + \frac{h^4}{24} f^{(4)}(\xi_+) \right) \right) \\ &= f''_i + \frac{h^2}{24} (f^{(4)}(\xi_-) + f^{(4)}(\xi_+)), \end{aligned}$$

где ξ_- , ξ_+ — некоторые точки из открытого интервала $]x_{i-1}, x_{i+1}[$. Поэтому если $f \in C^4[x_{i-1}, x_{i+1}]$, то справедлива оценка

$$|f''(x_i) - f_{x\bar{x},i}| \leq \frac{M_4}{12} h^2,$$

где $M_4 = \max_{\xi} |f^{(4)}(\xi)|$. Таким образом, порядок точности этой формулы равен 2 на функциях из C^4 .

Приведём ещё без вывода результат о погрешности формулы для вычисления второй производной вблизи края сетки (таблицы):

$$f''(x_i) = \frac{1}{h^2} (2f_i - 5f_{i+1} + 4f_{i+2} - f_{i+3}) + O(h^2),$$

$$f''(x_i) = \frac{1}{h^2} (f_{i-3} - 4f_{i-2} + 5f_{i-1} - 2f_i) + O(h^2).$$

Порядок этих формул всего лишь второй, откуда видна роль симметричности шаблона в трёхточечной формуле (2.65) с тем же порядком точности.

Что произойдёт, если дифференцируемая функция не будет иметь достаточную гладкость? Тогда мы не сможем выписывать необходимое количество членов разложения по формуле Тейлора, и потому полученный порядок точности формул с помощью метода разложений установить не сможем. Тот факт, что в этих условиях реальный порядок точности может быть в самом деле меньшим, чем для функций с высокой гладкостью, показывает следующий

Пример 2.8.2 Рассмотрим функцию $g(x) = x|x|$, которую эквивалентным образом можно задать в виде

$$g(x) = \begin{cases} x^2, & \text{если } x \geq 0, \\ -x^2, & \text{если } x \leq 0. \end{cases}$$

Её график изображён на Рис. 2.14.

Функция $g(x)$ дифференцируема всюду на числовой оси. При $x \neq 0$ она имеет производную, равную

$$g'(x) = (x|x|)' = x'|x| + x|x|' = |x| + x \operatorname{sgn} x = 2|x|,$$

а в нуле

$$g'(0) = \lim_{x \rightarrow 0} \frac{x|x|}{x} = 0.$$

Таким образом, производная $g'(x) = 2|x|$ всюду непрерывна. Но она недифференцируема в нуле, так что вторая производная $g''(0)$ уже не существует. Как следствие, $g(x) \in C^1$, но $g(x) \notin C^2$ на любом интервале, содержащем нуль.

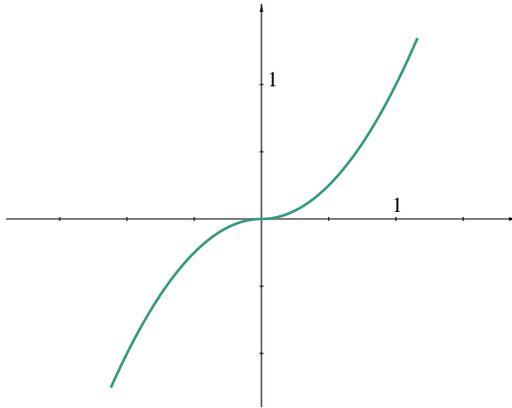


Рис. 2.14. График функции $y = x|x|$: увидеть разрыв её второй производной в нуле почти невозможно.

Воспользуемся для численного нахождения производной $g'(0)$ формулой центральной разности (2.64) на шаблоне с шагом h , симметричном относительно нуля:

$$g'(0) \approx \frac{g(h) - g(-h)}{2h} = \frac{h|h| - (-h)|-h|}{2h} = \frac{h^2 + h^2}{2h} = h.$$

Таким образом, при $h \rightarrow 0$ приближённое числовое значение производной стремится к $g'(0) = 0$ с первым порядком по h , а не вторым, как мы установили это ранее для дважды гладких функций. ■

2.8в Метод неопределённых коэффициентов

Метод неопределённых коэффициентов — это другой подход к получению формул численного дифференцирования, особенно удобный в многомерном случае, когда построение интерполяционного полинома становится непростым.

Предположим, что задан шаблон из $p + 1$ штук точек x_0, x_1, \dots, x_p . Станем искать приближённое выражение для производной от функции в виде линейной формы от значений функции, т. е. как

$$f^{(k)}(x) \approx \sum_{i=0}^p c_i f(x_i). \quad (2.75)$$

Она мотивируется тем обстоятельством, что дифференцирование любого порядка является операцией, линейной по значениям функции. Линейными формами от значений функции были, в частности, все полученные ранее формулы численного дифференцирования, начиная с (2.61) и кончая (2.72).

Коэффициенты c_i линейной формы постараемся подобрать так, чтобы эта формула являлась точной формулой для какого-то «достаточно представительного» набора функций. Например, в качестве таких «пробных функций» можно взять все полиномы степени не выше заданной, либо тригонометрические полиномы (2.4) какого-то фиксированного порядка и т. п. Рассмотрим ниже подробно случай алгебраических полиномов.

Возьмём $f(x)$ равной последовательным степеням переменной x , т. е. $1, x, x^2, \dots, x^q$ для некоторого фиксированного q . Если формула (2.75) обращается в точное равенство на этих «пробных функциях», то с учётом её линейности можно утверждать, что она будет точной для любого алгебраического полинома степени не выше q .

Каждое условие, выписанное для какой-то определённой степени x^j , $j = 0, 1, \dots, q$, является линейным соотношением для неизвестных коэффициентов c_i , и в целом мы приходим к системе линейных уравнений относительно c_i , $i = 0, 1, \dots, p$. Для разрешимости этой системы естественно взять число неизвестных равным числу уравнений, т. е. $q = p$.

Получающаяся система линейных уравнений имеет вид

$$\left\{ \begin{array}{l} c_0 + c_1 + \dots + c_p = 0, \\ c_0 x_0 + c_1 x_1 + \dots + c_p x_p = 0, \\ \vdots \quad \quad \quad \ddots \quad \quad \quad \vdots \quad \quad \quad \vdots \\ c_0 x_0^{k-1} + c_1 x_1^{k-1} + \dots + c_p x_p^{k-1} = 0, \\ c_0 x_0^k + c_1 x_1^k + \dots + c_p x_p^k = k!, \\ c_0 x_0^{k+1} + c_1 x_1^{k+1} + \dots + c_p x_p^{k+1} = (k+1)!x, \\ \vdots \quad \quad \quad \ddots \quad \quad \quad \vdots \quad \quad \quad \vdots \\ c_0 x_0^p + c_1 x_1^p + \dots + c_p x_p^p = p(p-1) \cdots (p-k+1) x^{p-k}. \end{array} \right. \quad (2.76)$$

В правых частях этой системы стоят k -е производные от $1, x, x^2, \dots, x^q$, а матрицей системы является матрица Вандермонда вида (2.7), которая неособенна для попарно различных узлов x_0, x_1, \dots, x_p . При этом система линейных уравнений однозначно разрешима относительно c_0, c_1, \dots, c_p для любой правой части, но содержательным является лишь случай $k \leq p$. В противном случае, если $k > p$, правая часть системы (2.76) оказывается нулевой, и, как следствие, система также имеет только бессодержательное нулевое решение. Этот факт имеет интуитивно ясное объяснение: нельзя построить формулу для вычисления производной k -го порядка от функции, используя значения этой функции не более чем в k точках.

Матрицы Вандермонда в общем случае являются плохообусловленными (см. §3.56). Но на практике решение системы (2.76) — вручную или на компьютере — обычно не приводит к большим ошибкам, так как порядок системы (2.76), равный порядку производной, бывает, как правило, небольшим.¹²

Пример 2.8.3 Построим формулу численного дифференцирования ■

Интересен вопрос о взаимоотношении метода неопределённых коэффициентов и рассмотренного ранее в §2.8а интерполяционного подхода

¹²На стр. 87 мы уже обсуждали вопрос о том, каков наивысший порядок производных, всё ещё имеющих содержательный смысл.

к численному дифференцированию. К примеру, Ш.Е. Микеладзе в книге [53] утверждает, что любая формула численного дифференцирования, полученная методом неопределённых коэффициентов, может быть выведена также с помощью интерполяционного подхода, отказывая методу неопределённых коэффициентов в оригинальности. Но нельзя отрицать также, что метод неопределённых коэффициентов конструктивно проще и технологичнее в применении, и уже только это обстоятельство оправдывает его существование.

2.8г Полная вычислительная погрешность численного дифференцирования

Рассмотрим поведение полной погрешности численного дифференцирования при расчётах на реальных вычислительных устройствах. Под *полной погрешностью* мы понимаем суммарную ошибку численного нахождения производной, вызванную как приближённым характером самого метода, так и неточностями вычислений на цифровых ЭВМ из-за неизбежных ошибок округления и т. п.

Предположим, к примеру, что первая производная функции вычисляется по формуле «разность вперёд»

$$f'(x_i) \approx f_{x,i} = \frac{f_{i+1} - f_i}{h}.$$

Как мы уже знаем, её погрешность

$$|f_{x,i} - f'(x_i)| \leq \frac{M_2 h}{2},$$

где $M_2 = \max_{\xi \in [a,b]} |f''(\xi)|$. Если значения функции вычисляются с ошибками, то вместо точных f_i и f_{i+1} мы получаем их приближённые значения \tilde{f}_i и \tilde{f}_{i+1} , такие что

$$|f_i - \tilde{f}_i| \leq \delta \quad \text{и} \quad |f_{i+1} - \tilde{f}_{i+1}| \leq \delta,$$

где через δ обозначена предельная абсолютная погрешность вычисления значений функции. Тогда в качестве приближённого значения производной мы должны взять

$$f'(x_i) \approx \frac{\tilde{f}_{i+1} - \tilde{f}_i}{h},$$

а предельную полную вычислительную погрешность $E(h, \delta)$ нахождения первой производной функции можно оценить следующим образом:

$$\begin{aligned}
 E(h, \delta) &= \left| \frac{\tilde{f}_{i+1} - \tilde{f}_i}{h} - f'(x_i) \right| \\
 &\leq \left| \frac{\tilde{f}_{i+1} - \tilde{f}_i}{h} - \frac{f_{i+1} - f_i}{h} \right| + \left| \frac{f_{i+1} - f_i}{h} - f'(x_i) \right| \\
 &\leq \left| \frac{(\tilde{f}_{i+1} - f_{i+1}) + (f_i - \tilde{f}_i)}{h} \right| + \frac{M_2 h}{2} \\
 &\leq \frac{|f_{i+1} - \tilde{f}_{i+1}| + |f_i - \tilde{f}_i|}{h} + \frac{M_2 h}{2} = \frac{2\delta}{h} + \frac{M_2 h}{2}.
 \end{aligned} \tag{2.77}$$

Отметим, во-первых, что эта оценка, достижима при подходящем сочетании знаков фигурирующих в неравенствах величин, коль скоро достижимо используемое в преобразованиях неравенство треугольника $|a + b| \leq |a| + |b|$ и достижима оценка погрешности (2.74) для формулы «разность вперёд». Во-вторых, оценка не стремится к нулю при уменьшении шага h , так как первое слагаемое неограниченно увеличивается при $h \rightarrow 0$. В целом, функция $E(h, \delta)$ при фиксированном δ имеет минимум, определяемый условием

$$\frac{\partial E(h, \delta)}{\partial h} = \frac{\partial}{\partial h} \left(\frac{2\delta}{h} + \frac{M_2 h}{2} \right) = -\frac{2\delta}{h^2} + \frac{M_2}{2} = 0.$$

То есть, оптимальное значение шага численного дифференцирования, при котором достигается минимальная полная погрешность, равно

$$h^* = 2\sqrt{\delta/M_2}, \tag{2.78}$$

и брать меньший шаг численного дифференцирования смысла нет. Само значение достигаемой при этом полной погрешности есть $E(h^*, \delta) = 2\sqrt{\delta M_2}$.

Пример 2.8.4 Пусть в арифметике двойной точности с плавающей точкой, реализованной согласно стандарту IEEE 754/854, численно находится производная функции, выражение для которой требует десять

вычислений, а модуль второй производной ограничен сверху величиной $M_2 = 10$. Погрешность отдельной арифметической операции можно считать приближённо равной половине расстояния между соседними машинно представимыми числами, т.е. примерно 10^{-16} в районе единицы. Наконец, пусть абсолютная погрешность вычисления функции складывается из сумм абсолютных погрешностей каждой операции, так что $\delta \approx 10 \cdot 10^{-16} = 10^{-15}$ при аргументах порядка единицы.

Тогда в соответствии с формулой (2.78) имеем $h^* = 2\sqrt{\delta/M_2} = 2 \cdot 10^{-8}$, т.е. брать шаг сетки меньше 10^{-8} смысла не имеет. ■

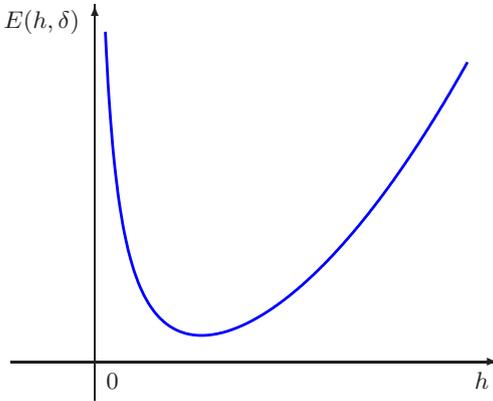


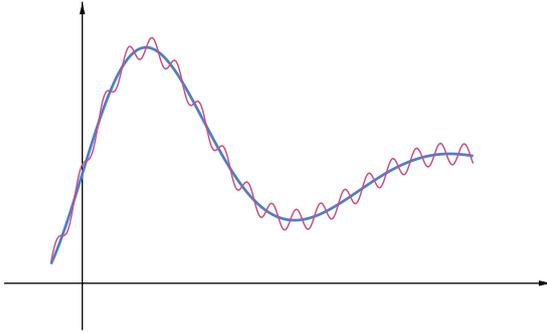
Рис. 2.15. Типичный график полной погрешности численного дифференцирования

Совершенно аналогичная ситуация имеет место и при использовании других формул численного дифференцирования. Производная k -го порядка на равномерной сетке шага h определяется в общем случае формулой вида¹³

$$f^{(k)}(x) = h^{-k} \sum_i c_i f(x_i) + R_k(f, x), \quad (2.79)$$

где $c_i = O(1)$ при $h \rightarrow 0$. Если эта формула имеет порядок точности p , то её остаточный член оценивается как $R_k(f, x) \approx c(x)h^p$. Этот остаточный член определяет «идеальную» погрешность численного диффе-

¹³Для примера можно взглянуть на те формулы, которые приведены в §2.8а.

Рис. 2.16. Возмущение функции добавкой $\frac{1}{n} \sin(nx)$.

ренцирования в отсутствие ошибок вычисления функции, и он неограниченно убывает при $h \rightarrow 0$.

Но если погрешность вычисления значений функции $f(x_i)$ в узлах равна δ , то в правой части (2.79) возникает ещё член, абсолютная величина которого совершенно аналогично (2.77) оценивается сверху как

$$\delta h^{-k} \sum_i |c_i|.$$

Она неограниченно возрастает при $h \rightarrow 0$. В целом график полной вычислительной погрешности численного дифференцирования выглядит в этом случае примерно так, как на Рис. 2.15.

Практический вывод из сказанного состоит в том, что существует оптимальный шаг h численного дифференцирования, минимизирующий полную вычислительную погрешность, и брать слишком маленькое значение шага h в практических расчётах нецелесообразно.

Потенциально сколь угодно большое возрастание погрешности численного дифференцирования, в действительности, является отражением более глубокого факта *некорректности* задачи дифференцирования (см. §1.3). Её решение не зависит непрерывно от входных данных, и это демонстрируют простые примеры. Если $f(x)$ — исходная функция, производную которой нам требуется найти, то возмущённая функция $f(x) + \frac{1}{n} \sin(nx)$ при $n \rightarrow \infty$ будет равномерно сходиться к исходной, тогда как её производная

$$f'(x) + \cos(nx)$$

не сходится к производной $f'(x)$ (см. Рис. 2.16). При возмущении исходной функции слагаемым $\frac{1}{n} \sin(n^2x)$ производная вообще может сколь угодно сильно отличаться от производной исходной функции.

2.9 Алгоритмическое дифференцирование

Пусть $u = u(x)$ и $v = v(x)$ — некоторые выражения от переменной x , из которых далее с помощью сложения, вычитания, умножения или деления конструируется более сложное выражение. Напомним правила дифференцирования выражений, образованных с помощью элементарных арифметических операций:

$$(u + v)' = u' + v', \quad (2.80)$$

$$(u - v)' = u' - v', \quad (2.81)$$

$$(uv)' = u'v + uv', \quad (2.82)$$

$$\left(\frac{u}{v}\right)' = \frac{u'v - uv'}{v^2}. \quad (2.83)$$

Из них следует, что численное значение производной для сложного выражения мы можем найти, зная лишь значения образующих его подвыражений и их производных.

Сделанное наблюдение подсказывает идею ввести на множестве пар вида (u, u') , которые составлены из значений выражения и его производной, арифметические операции по правилам, следующим из формул (2.80)–(2.83):

$$(u, u') + (v, v') = (u + v, u' + v'), \quad (2.84)$$

$$(u, u') - (v, v') = (u - v, u' - v'), \quad (2.85)$$

$$(u, u') \cdot (v, v') = (uv, u'v + uv'), \quad (2.86)$$

$$\frac{(u, u')}{(v, v')} = \left(\frac{u}{v}, \frac{u'v - uv'}{v^2}\right). \quad (2.87)$$

Первые члены пар преобразуются просто в соответствии с применяемой арифметической операцией, а операции над вторыми членами пар — это в точности копии правил (2.80)–(2.83). Если для заданного выражения мы начнём вычисления по выписанным формулам (2.84)–(2.87), заменив исходную переменную x на пары $(x, 1)$, а константы c — на

пары вида $(c, 0)$, то на выходе получим пару, состоящую из численных значений выражения и производной от него в точке x .

Это рассуждение очевидно обобщается на случай, когда функция зависит от нескольких переменных.

Помимо арифметических операций интересующее нас выражение может содержать вхождения элементарных функций. Для них в соответствии с формулами дифференциального исчисления можем определить действия над парами следующим образом

$$\begin{aligned}\exp((u, u')) &= (\exp u, u' \exp u), \\ \sin((u, u')) &= (\sin u, u' \cos u), \\ ((u, u'))^2 &= (u^2, 2uu'), \\ ((u, u'))^3 &= (u^3, 3u^2u') \text{ и т.д.}\end{aligned}$$

Арифметику пар вида (u, u') с операциями (2.84)–(2.87) называют *дифференциальной арифметикой*, а основанный на её использовании способ вычисления значений производных носит название *алгоритмического дифференцирования*. Нередко используют также термин «автоматическое дифференцирование».

Строго говоря, мы рассмотрели один из возможных способов организации алгоритмического дифференцирования, который называют *прямым режимом*. Существует ещё и *обратный режим* алгоритмического дифференцирования.

Описанную выше идею можно применить к вычислению вторых производных. Но теперь вместо дифференциальной арифметики пар чисел (u, u') нам необходимо будет оперировать с числовыми тройками вида (u, u', u'') , поскольку в формулах для вторых производных функции фигурируют значения самой функции и её первых и вторых производных.

Идея алгоритмического дифференцирования может быть распространена на вычисление разделённых разностей (наклонов) функций, а также на вычисление интервальных расширений производных и наклонов (см., к примеру, [65]).

2.10 Приближение функций

2.10а Обсуждение постановки задачи

В этом параграфе мы займёмся задачей приближения функций. К ней естественно приходят в ситуациях, где методы интерполирования по различным причинам не удовлетворяют практику. Эти причины могут носить чисто технический характер. К примеру, гладкость сплайна может оказаться недостаточной, либо его построение — слишком сложным. Степень обычного интерполяционного полинома может быть неприемлемо высокой для данного набора узлов интерполяции (а высокая степень — это трудности при вычислении значений полинома и его большая изменчивость). Но причины отказа от интерполяции могут иметь также принципиальный характер. В частности, это происходит, если значения функции в узлах x_0, x_1, \dots, x_n известны неточно. В этих условиях целесообразна коррекция самой постановки задачи.

Именно, имеет смысл отказаться от требования, чтобы восстанавливаемая функция g была точно равна значениям f_i в узлах x_0, x_1, \dots, x_n , допустив, к примеру, для g принадлежности её значений некоторым интервалам, т. е. $g(x_i) \in [\underline{f}_i, \overline{f}_i]$, $i = 0, 1, \dots, n$, $\underline{f}_i \leq \overline{f}_i$. Наглядно-геометрически это означает построение функции $g(x)$ из заданного класса \mathcal{G} , которая в каждом узле сетки x_i , $i = 0, 1, \dots, n$, проходит через некоторый «коридор» $[\underline{f}_i, \overline{f}_i]$, см. Рис. 2.17.

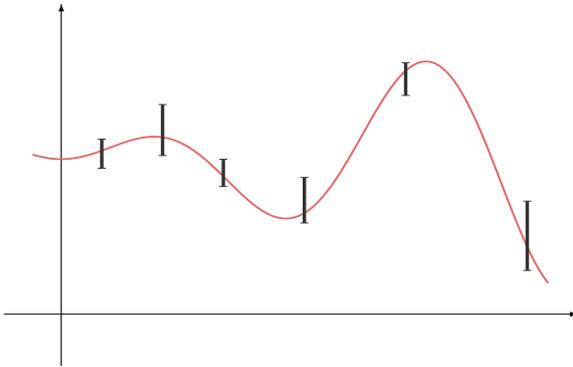


Рис. 2.17. Интерполяция функции, заданной с погрешностью

Более общая постановка этой задачи предусматривает наличие неко-

торой метрики (расстояния), которую мы будем обозначать через dist , и с помощью которой можно измерять отклонение вектора значений $(g(x_0), g(x_1), \dots, g(x_n))^T$ функции $g(x)$ в узлах сетки от вектора заданных значений $(f_0, f_1, \dots, f_n)^T$. Напомним, что на множестве Y , образованном элементами произвольной природы, *расстоянием* (или *метрикой*) называется определенная на декартовом произведении $Y \times Y$ функция dist с неотрицательными вещественными значениями, удовлетворяющая для любых $f, g, h \in Y$ следующим условиям:

- (1) $\text{dist}(f, g) = 0$ тогда и только тогда, когда $f = g$,
- (2) $\text{dist}(f, g) = \text{dist}(g, f)$ — симметричность,
- (3) $\text{dist}(f, h) \leq \text{dist}(f, g) + \text{dist}(g, h)$ — неравенство треугольника.

Фактически, в рассмотренной выше ситуации dist — это какое-то расстояние на пространстве \mathbb{R}^{n+1} всех $(n+1)$ -мерных вещественных векторов. Соответствующая постановка *задачи приближения* (аппроксимации) будет звучать тогда так:

Для заданного набора узлов $x_i, i = 0, 1, \dots, n$, на интервале $[a, b]$ и соответствующих им значений $f_i, i = 0, 1, \dots, n$, и $\epsilon > 0$, найти функцию $g(x)$ из класса функций \mathcal{G} , такую что $\text{dist}(f, g) < \epsilon$, где $f = (f_0, f_1, \dots, f_n)^T$ и $g = (g(x_0), g(x_1), \dots, g(x_n))^T$.

При этом $g(x)$ называют *приближающей* (аппроксимирующей) функцией, Важнейшей модификацией поставленной задачи служит *задача наилучшего приближения*, когда величина ϵ не фиксируется и ищут приближающую (аппроксимирующую) функцию $g(x)$, которая доставляет минимум расстоянию $\text{dist}(f, g)$.

Согласно классификации §2.1, выписанные выше формулировки являются дискретными вариантами общей задачи о приближении функции, в которой дискретный набор узлов x_0, x_1, \dots, x_n уже не фигурирует, а отклонение одной функции от другой измеряется «на всей» области их определения:

Для заданных $\epsilon > 0$, функции $f(x)$ из \mathcal{F} и метрики dist найти функцию $g(x)$ из класса функций \mathcal{G} , такую что $\text{dist}(f, g) < \epsilon$.

Соответствующая общая формулировка задачи о наилучшем прибли-

жении ставится так:

Для заданных функции $f(x)$ из класса функций \mathcal{F} и метрики dist найти функцию $g(x)$ из класса \mathcal{G} , на которой достигается нижняя грань расстояний от $f(x)$ до функций из \mathcal{G} , т. е. удовлетворяющую условию $\text{dist}(f, g) = \inf_{h \in \mathcal{G}} \text{dist}(f, h)$. (2.88)

Решение g этой задачи, если оно существует, называется *наилучшим приближением* для f в классе \mathcal{G} . Отметим, что в каждом конкретном случае существование элемента наилучшего приближения требует отдельного исследования.

Отметим, что задачу приближения функций, значения которых заданы приближённо, часто называют (особенно в практических приложениях) *задачей сглаживания*, поскольку получаемая приближающая функция, как правило, действительно «сглаживает» выбросы данных, вызванные случайными ошибками и т. п.

До сих пор ничего не было сказано о выборе классов функций \mathcal{F} и \mathcal{G} , и в наших формулировках они могут быть весьма произвольными. Но чаще всего предполагают, что $\mathcal{F} \supseteq \mathcal{G}$, и, кроме того, наделяют \mathcal{F} и \mathcal{G} структурой линейного пространства с некоторой нормой $\|\cdot\|$. Именно в ней измеряют отклонение функций (непрерывного или дискретного аргумента) друг от друга, так что

$$\text{dist}(f, g) = \|f - g\|.$$

Соответственно, в задаче наилучшего приближения функции f ищется такой элемент $g \in \mathcal{G}$, на котором достигается $\inf_{h \in \mathcal{G}} \|f - h\|$.

Рассмотренные выше постановки задач дают начало большим и важным разделам математики, в совокупности образующим теорию приближения функций (называемую также теорией аппроксимации). Её ветвью является, в частности, теория равномерного приближения, когда отклонение функций оценивается в норме $\|f\| = \max_{x \in [a, b]} |f(x)|$ (см. [45, 59]). Выбор различных норм (т. е. различных мер отклонения функций друг от друга) и различных классов функций обуславливает огромное разнообразие задач теории приближения.

2.106 Существование и единственность решения задачи приближения

Некоторые свойства решения задачи наилучшего приближения функций можно вывести уже из абстрактной формулировки. В частности, это касается существования решения, а также единственности решения при некоторых дополнительных условиях на норму.

Предложение 2.10.1 Пусть X — нормированное линейное пространство, а U — его конечномерное линейное подпространство. Тогда для любого $f \in X$ существует элемент наилучшего приближения $u \in U$.

Доказательство. Пусть размерность U равна m . Зафиксировав некоторый базис $\phi_1, \phi_2, \dots, \phi_m$ подпространства U , введём функцию

$$r(a_1, a_2, \dots, a_m) = \left\| f - \sum_{j=1}^m a_j \phi_j \right\|.$$

Предложение будет доказано, если мы обоснуем тот факт, что функция $r : \mathbb{R}^m \rightarrow \mathbb{R}$ достигает своего наименьшего значения на \mathbb{R}^m .

Прежде всего покажем, что функция r непрерывно зависит от своих аргументов:

$$\begin{aligned} & \left| r(b_1, b_2, \dots, b_m) - r(a_1, a_2, \dots, a_m) \right| \\ & \leq \left\| \left\| f - \sum_{j=1}^m b_j \phi_j \right\| - \left\| f - \sum_{j=1}^m a_j \phi_j \right\| \right\| \\ & \leq \left\| \sum_{j=1}^m (b_j - a_j) \phi_j \right\| \leq \sum_{j=1}^m |b_j - a_j| \|\phi_j\| \\ & \leq \max_{1 \leq j \leq m} |b_j - a_j| \cdot \sum_{j=1}^m \|\phi_j\|. \end{aligned}$$

Следовательно, при $b_j \rightarrow a_j$ разность между $r(b_1, b_2, \dots, b_m)$ и $r(a_1, a_2, \dots, a_m)$ также будет стремиться к нулю.

Следующим шагом доказательства продемонстрируем, что функция $r(x)$ может достигать своего минимума лишь на некотором подмножестве всего пространстве \mathbb{R}^m , которое к тому же компактно.



Элемент наилучшего приближения, вообще говоря, может быть неединственным. Но при определённых условиях мы можем гарантировать его единственность, опираясь лишь на свойства пространства X .

Нормированное пространство X с нормой $\|\cdot\|$ называют *строго нормированным*, если для произвольных $x, y \in X$ из равенства $\|x + y\| = \|x\| + \|y\|$ следует существование такого скаляра $\alpha \in \mathbb{R}$, что $y = \alpha x$. Иными словами, в таком пространстве равенство в неравенстве треугольника возможно лишь для коллинеарных векторов.

Пример 2.10.1 Строго нормированным пространством является \mathbb{R}^2 с евклидовой нормой $\|x\|_2 = \sqrt{x_1^2 + x_2^2}$. Но нормы $\|x\|_1 = |x_1| + |x_2|$ и $\|x\|_\infty = \max\{|x_1|, |x_2|\}$ на \mathbb{R}^2 , которые эквивалентны норме $\|\cdot\|_2$ (см. §3.3б), не порождают строго нормированное пространство. ■

Предложение 2.10.2 Пусть X — строго нормированное линейное пространство, а U — его конечномерное линейное подпространство. Тогда для любого $f \in X$ элемент его наилучшего приближения $u \in U$ единствен.

Доказательство. Предположим, что для элемента f существуют два наилучших приближения

$$u' = \sum_{i=1}^m u'_i \phi_i \quad \text{и} \quad u'' = \sum_{i=1}^m u''_i \phi_i, \quad (2.89)$$

которые определяются наборами коэффициентов $(u'_1, u'_2, \dots, u'_m)$ и $(u''_1, u''_2, \dots, u''_m)$ разложения по базису $\phi_i, i = 1, 2, \dots, m$. При этом

$$\left\| f - \sum_{i=1}^m u'_i \phi_i \right\| = \left\| f - \sum_{i=1}^m u''_i \phi_i \right\| = \mu \geq 0,$$

где μ — величина наименьшего отклонения f от u' и u'' .

Возьмём середину отрезка прямой, соединяющей u' и u'' , т. е. точку, у которой компоненты разложения по векторам ϕ_i , равны $\frac{1}{2}(u'_i + u''_i)$.

Имеем

$$\begin{aligned} \left\| f - \sum_{i=1}^m \frac{1}{2}(u'_i + u''_i) \phi_i \right\| &= \left\| \frac{1}{2} \left(f - \sum_{i=1}^m u'_i \phi_i \right) + \frac{1}{2} \left(f - \sum_{i=1}^m u''_i \phi_i \right) \right\| \\ &\leq \frac{1}{2} \left\| f - \sum_{i=1}^m u'_i \phi_i \right\| + \frac{1}{2} \left\| f - \sum_{i=1}^m u''_i \phi_i \right\| = \mu. \end{aligned}$$

Строгого неравенства здесь быть не может, что очевидно для $\mu = 0$, а для $\mu > 0$ означало бы существование элемента, приближающего f лучше, чем два элемента наилучшего приближения u' и u'' . Поэтому необходимо должно выполняться равенство

$$\begin{aligned} \left\| \frac{1}{2} \left(f - \sum_{i=1}^m u'_i \phi_i \right) + \frac{1}{2} \left(f - \sum_{i=1}^m u''_i \phi_i \right) \right\| \\ = \frac{1}{2} \left\| f - \sum_{i=1}^m u'_i \phi_i \right\| + \frac{1}{2} \left\| f - \sum_{i=1}^m u''_i \phi_i \right\|. \end{aligned}$$

Но если рассматриваемое пространство — строго нормированное, то из полученного равенства следует

$$f - \sum_{i=1}^m u'_i \phi_i = \alpha \left(f - \sum_{i=1}^m u''_i \phi_i \right) \quad (2.90)$$

для некоторого вещественного α .

В случае, когда $\alpha \neq 1$, получаем

$$f = \frac{1}{1-\alpha} \cdot \sum_{i=1}^m (u'_i - \alpha u''_i) \phi_i,$$

т. е. f точно представляется в виде линейной комбинации базисных векторов ϕ_i . Тогда $\mu = 0$, и в силу единственности разложения по базису должно быть $u'_i = u''_i$. Следовательно, для f в действительности существует всего лишь одно наилучшее приближение.

В остающемся случае $\alpha = 1$ для выполнения равенства (2.90) необходимо $u'_i = u''_i$, и тогда два элемента наилучшего приближения (2.89) также должны совпадать. ■

2.10в Задача приближения в евклидовом пространстве

Рассмотрим подробно важный частный случай задачи о наилучшем приближении (2.88), в котором

- класс \mathcal{F} — линейное нормированное пространство функций, на котором задано скалярное произведение $\langle \cdot, \cdot \rangle$, и с его помощью норма в \mathcal{F} определяется как $\|f\| = \sqrt{\langle f, f \rangle}$,
- класс функций $\mathcal{G} \subseteq \mathcal{F}$, из которого выбирается искомый элемент наилучшего приближения, является конечномерным подпространством в \mathcal{F} .

Напомним, что конечномерные линейные векторные пространства, в которых определено скалярное произведение, называются *евклидовыми пространствами*. Бесконечномерные линейные векторные пространства со скалярным произведением называются *гильбертовыми пространствами* при дополнительном условии полноты, т. е. существования в них предела всякой фундаментальной последовательности относительно нормы, порождённой этим скалярным произведением. Гильбертовы пространства являются ближайшим обобщением пространств с привычной нам геометрией.

В условиях постановки задачи, описанной в начале раздела, будем предполагать, что известен $\{\varphi_j\}_{j=1}^m$ — базис m -мерного линейного подпространства $\mathcal{G} \subseteq \mathcal{F}$. Мы ищем приближение g для элемента $f \in \mathcal{F}$ в виде

$$g = \sum_{j=1}^m c_j \varphi_j. \quad (2.91)$$

где $c_j, j = 1, 2, \dots, m$ — неизвестные коэффициенты, подлежащие определению. Если через Φ обозначить квадрат нормы отклонения f от g , то имеем

$$\begin{aligned} \Phi &= \|f - g\|^2 = \langle f - g, f - g \rangle \\ &= \langle f, f \rangle - 2\langle f, g \rangle + \langle g, g \rangle \\ &= \langle f, f \rangle - 2 \sum_{j=1}^m c_j \langle f, \varphi_j \rangle + \sum_{j=1}^m \sum_{k=1}^m c_j c_k \langle \varphi_j, \varphi_k \rangle. \end{aligned} \quad (2.92)$$

Как видим, Φ есть квадратичная форма от аргументов c_1, c_2, \dots, c_m плюс ещё некоторые линейные члены относительно c_j и постоянное слагаемое $\langle f, f \rangle$. Её особенностью является то обстоятельство, что для всех значений аргументов функция Φ принимает только неотрицательные значения. Покажем, что она достигает своего минимума.

Пусть $\Pi_R = \{x \in \mathbb{R}^m \mid \sqrt{x_1^2 + x_2^2 + \dots + x_m^2} \leq R\}$ — замкнутый шар радиуса R с центром в нуле относительно евклидова расстояния. Рассмотрим поведение $\min_{c \in \Pi_R} \Phi(c)$ в зависимости от R . При увеличении R значение этого минимума не возрастает, но, в действительности, оно не может уменьшаться, начиная с некоторого R .

В самом деле, после приведения к «главным осям» квадратичная форма в составе Φ обязательно должна получить вид суммы квадратов с положительными коэффициентами, так как иначе вся Φ была бы неограниченной снизу. Но сумма квадратов неограниченно возрастает при увеличении расстояния аргумента $c = (c_0, c_1, \dots, c_m)$ до нуля, причём растёт быстрее линейных членов. Следовательно, при достаточно больших значениях R его увеличение уже не окажет никакого влияния на $\min_{c \in \Pi_R} \Phi(c)$, и потому мы сможем утверждать, что $\min_{c \in \mathbb{R}^m} \Phi(c)$ достигается в некотором шаре Π_R . В силу компактности множества Π_R это означает, что $\min \Phi(c)$ действительно достигается в некоторой конечной точке из \mathbb{R}^m .

Для определения минимума функции Φ продифференцируем её по c_j , $j = 1, 2, \dots, m$, и приравняем полученные производные к нулю:

$$\frac{\partial \Phi}{\partial c_j} = -2\langle f, \varphi_j \rangle + 2 \sum_{k=1}^m c_k \langle \varphi_j, \varphi_k \rangle = 0. \quad (2.93)$$

Множитель 2 при сумме всех $c_k \langle \varphi_j, \varphi_k \rangle$ появляется оттого, что в двойной сумме из выражения (2.92) слагаемое с c_j возникает дважды: один раз с коэффициентом $\langle \varphi_j, \varphi_k \rangle$, а другой раз — с коэффициентом $\langle \varphi_k, \varphi_j \rangle$.

В целом, из равенств (2.93) для определения c_j получаем систему линейных алгебраических уравнений

$$\sum_{k=1}^m \langle \varphi_j, \varphi_k \rangle c_k = \langle f, \varphi_j \rangle, \quad j = 1, 2, \dots, m. \quad (2.94)$$

Матрица её коэффициентов

$$\Gamma(\varphi_1, \varphi_2, \dots, \varphi_m) = \begin{pmatrix} \langle \varphi_1, \varphi_1 \rangle & \langle \varphi_1, \varphi_2 \rangle & \dots & \langle \varphi_1, \varphi_m \rangle \\ \langle \varphi_2, \varphi_1 \rangle & \langle \varphi_2, \varphi_2 \rangle & \dots & \langle \varphi_2, \varphi_m \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \varphi_m, \varphi_1 \rangle & \langle \varphi_m, \varphi_2 \rangle & \dots & \langle \varphi_m, \varphi_m \rangle \end{pmatrix}, \quad (2.95)$$

называется, как известно, *матрицей Грама* системы векторов $\varphi_1, \varphi_2, \dots, \varphi_m$. Из курса линейной алгебры и аналитической геометрии читателю должно быть известно, что матрица Грама — это симметричная матрица, неособенная тогда и только тогда, когда векторы $\varphi_1, \varphi_2, \dots, \varphi_m$ линейно независимы (см., к примеру, [31]). При выполнении этого условия матрица Грама является ещё и положительно определённой. Таким образом, решение задачи наилучшего среднеквадратичного приближения существует и единственно, если $\varphi_1, \varphi_2, \dots, \varphi_m$ образуют базис в подпространстве \mathcal{G} .

Обратимся к практическим аспектам реализации развитого выше метода и обсудим свойства системы уравнений (2.94). Особенно интересна устойчивость её решения к возмущениям в данных и погрешностям вычислений на цифровых ЭВМ.

Наиболее простой вид матрица Грама имеет в случае, когда базисные функции φ_j ортогональны друг другу, т. е. когда $\langle \varphi_j, \varphi_k \rangle = 0$ при $j \neq k$. При этом система линейных уравнений (2.94) становится диагональной и решается тривиально. Соответствующее наилучшее приближение имеет тогда вид суммы

$$g = \sum_{j=1}^m c_j \varphi_j, \quad \text{где } c_j = \frac{\langle f, \varphi_j \rangle}{\langle \varphi_j, \varphi_j \rangle}, \quad j = 1, 2, \dots, m, \quad (2.96)$$

и, как известно, называется (конечным) *рядом Фурье* для f по ортогональной системе векторов $\{\varphi_j\}_{j=1}^m$. Коэффициенты c_j из (2.96) называют при этом *коэффициентами Фурье* разложения функции f .

Кроме того, в случае ортогонального и близкого к ортогональному базиса $\{\varphi_j\}_{j=1}^m$ решение системы (2.94) устойчиво к возмущениям в правой части и неизбежным погрешностям вычислений. Но в общем случае базис линейного подпространства \mathcal{G} может сильно отличаться от ортогонального, и тогда свойства системы уравнений (2.94) могут быть плохими в том смысле, что её решение будет чрезвычайно чувствительным к возмущениям и погрешностям.

2.10г Среднеквадратичное приближение функций

В этом разделе мы применим развитый выше общий подход к конкретной задаче наилучшего среднеквадратичного приближения функций, заданных на интервале вещественной оси.

Приближение функций в норме, порождённой скалярным произведением, часто называют *среднеквадратичным приближением* или просто *квадратичным*. Дело в том, что в конечномерной ситуации скалярным произведением векторов $f = (f_0, f_1, \dots, f_n)^\top$ и $g = (g_0, g_1, \dots, g_n)^\top$ обычно берут

$$\langle f, g \rangle = \frac{1}{n+1} \sum_{i=0}^n \varrho_i f_i g_i, \quad (2.97)$$

для некоторого положительного весового вектора $\varrho = (\varrho_0, \varrho_1, \dots, \varrho_n)^\top$, $\varrho_i > 0$. То есть, соответствующая норма $\|\cdot\|$ такова, что расстояние одной функции до другой есть

$$\text{dist}(f, g) = \|f - g\| = \left(\frac{1}{n+1} \sum_{i=0}^n \varrho_i (f_i - g_i)^2 \right)^{1/2}, \quad (2.98)$$

— усреднение квадратов разностей компонент с какими-то весовыми множителями ϱ_i , $i = 0, 1, \dots, n$. В частности, если известна информация о точности задания отдельных значений функции f_i , то веса ϱ_i можно назначать так, чтобы отразить величину этой точности, сопоставляя значениям f_i с большей точностью больший вес.

Нормирующий множитель $\frac{1}{n+1}$ при суммах в (2.97) и (2.98) удобно брать для того, чтобы с ростом размерности n (при росте количества наблюдений, измельчении сетки и т. п.) ограничить рост величины скалярного произведения и нормы, обеспечив тем самым соизмеримость результатов при различных n .

Если f и g — функции непрерывного аргумента, то обычно полагают скалярное произведение равным

$$\langle f, g \rangle = \int_a^b \varrho(x) f(x) g(x) dx, \quad (2.99)$$

для некоторой весовой функции $\varrho(x) > 0$. Это выражение с точностью до множителя можно рассматривать как предел выражения (2.97) при

$n \rightarrow \infty$, так как в (2.97) легко угадываются интегральные суммы Римана для интеграла (2.99) по интервалу $[a, b]$ единичной длины при его равномерном разбиении на подинтервалы. Тогда аналогом (2.98) является расстояние между функциями

$$\text{dist}(f, g) = \|f - g\| = \left(\int_a^b \varrho(x)(f(x) - g(x))^2 dx \right)^{1/2}. \quad (2.100)$$

В связи с решением рассматриваемой задачи приближения функций часто используют термин *метод наименьших квадратов*. Фактически, это общее название целого семейства идейно близких методов построения приближений, которые основаны на минимизации суммы квадратов отклонений компонент исходного вектора от приближающего. В случае, когда рассматривается приближение функций (или вообще элементов) каких-то абстрактных пространств, естественным обобщением минимизации суммы квадратов является нахождение минимума нормы, порождённой скалярным произведением.

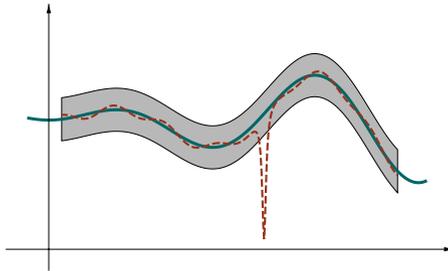


Рис. 2.18. Различие равномерного и интегрального (в частности, среднеквадратичного) отклонений функций друг от друга

Пример 2.10.2 В качестве примера практического возникновения задачи среднеквадратичного приближения рассмотрим тепловое действие тока $I(t)$ в проводнике сопротивлением R . Мгновенная тепловая мощность, как известно из теории электричества, равна при этом $I^2(t)R$, а полное количество теплоты, выделившееся между моментами времени a и b , равно

$$\int_a^b I^2(t)R dt.$$

Если мы хотим, скажем, минимизировать тепловыделение рассматриваемого участка электрической цепи, то нам нужно искать такой режим её работы, при котором достигался бы минимум выписанного интеграла, т. е. среднеквадратичного значения тока. В электротехнике его называют *действующим* или *эффективным* значением силы тока. ■

Линейным пространством \mathcal{F} , элементы которого мы будем приближать, выступит пространство всех функций, квадрат (т. е. степень 2) которых интегрируем на интервале $[a, b]$ с заданным весом $\varrho(x)$. Его называют пространством $\mathcal{L}^2[a, b]$, и нам сначала требуется показать, что оно в самом деле является линейным.

Ясно, что если $f \in \mathcal{L}^2[a, b]$, то для любого скаляра c функция $cf(x)$ также интегрируема с квадратом на $[a, b]$. Далее, с силу очевидного неравенства

$$2|f(x)g(x)| \leq (f(x))^2 + (g(x))^2$$

из интегрируемости второй степени функций $f(x)$ и $g(x)$ с весом $\varrho(x)$ следует интегрируемость их произведения на $[a, b]$. Поэтому существует интеграл

$$\int_a^b \varrho(x) (f(x) + g(x))^2 dx = \int_a^b \varrho(x) (f(x))^2 dx + 2 \int_a^b \varrho(x) f(x)g(x) dx + \int_a^b \varrho(x) (g(x))^2 dx,$$

т. е. сумма $f(x) + g(x)$ также имеет интегрируемый с весом $\varrho(x)$ квадрат. Это завершает доказательство линейности пространства $\mathcal{L}^2[a, b]$. Скалярное произведение в нём задаётся выражением (2.99). В курсах функционального анализа показывается, что если интегрирование понимается в смысле Лебега, то $\mathcal{L}^2[a, b]$ — гильбертово пространство, т. е. дополнительно обладает свойством полноты. По этой причине оно очень популярно в самых различных математических дисциплинах, от теории уравнений в частных производных до статистики.

Пример 2.10.3 Рассмотрим задачу о среднеквадратичном приближении функций из $\mathcal{L}^2[0, 1]$ с единичным весом полиномами фиксированной степени m . Скалярное произведение определяется как

$$\langle f, g \rangle = \int_0^1 f(x)g(x) dx,$$

а нормой берём

$$\|f\| = \int_0^1 (f(x))^2 dx.$$

Соответственно, расстояние между функциями определяется тогда как

$$\text{dist}(f, g) = \|f - g\| = \left(\int_0^1 (f(x) - g(x))^2 dx \right)^{1/2}.$$

Если в качестве базиса в линейном подпространстве полиномов мы возьмём последовательные степени

$$1, \quad x, \quad x^2, \quad \dots, \quad x^m,$$

то на месте (i, j) в матрице Грама (2.95) размера $(m+1) \times (m+1)$ будет стоять элемент

$$\int_0^1 x^{i-1} x^{j-1} dx = \frac{x^{i+j-1}}{i+j-1} \Big|_0^1 = \frac{1}{i+j-1}, \quad i, j = 1, 2, \dots, m+1$$

(сдвиг показателей степени на (-1) вызван тем, что строки и столбцы матрицы нумеруются, начиная с единицы, а не с нуля, как последовательность степеней). Матрица $H = (h_{ij})$ с элементами $h_{ij} = 1/(i+j-1)$, имеющая вид

$$\begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \dots & \frac{1}{m+1} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \dots & \frac{1}{m+2} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \dots & \frac{1}{m+3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{m+1} & \frac{1}{m+2} & \frac{1}{m+3} & \dots & \frac{1}{2m+1} \end{pmatrix},$$

называется *матрицей Гильберта*, и она является исключительно плохо обусловленной матрицей (см. §3.5б). Иными словами, решение СЛАУ с этой матрицей является непростой задачей, которая очень чувствительна к влиянию погрешностей в данных и вычислениях. ■

Пример 2.10.4 Пусть k и l — натуральные числа. Поскольку

$$\int_0^{2\pi} \sin(kx) \cos(lx) dx = 0,$$

для любых k, l , и

$$\int_0^{2\pi} \sin(kx) \sin(lx) dx = 0, \quad \int_0^{2\pi} \cos(kx) \cos(lx) dx = 0,$$

для $k \neq l$, то базис из тригонометрических полиномов вида

$$1, \quad \cos(2\pi kx), \quad \sin(2\pi kx), \quad k = 1, 2, \dots,$$

является ортогональным на $[0, 1]$ относительно скалярного произведения (2.99) с весом $\varrho(x) = 1$. Иными словами, этот базис очень хорош в вычислительном отношении для построения среднеквадратичных приближений. ■

Более детальный теоретический анализ и практический опыт показывают, что в методе наименьших квадратов в качестве базиса $\varphi_1, \varphi_2, \dots, \varphi_n$ линейного подпространства $\mathcal{G} \subset \mathcal{F}$ имеет смысл брать системы элементов, ортогональных по отношению к какому-то скалярному произведению (возможно, другому), так как это служит гарантией «разумной малости» внедиагональных элементов матрицы Грама и, как следствие, её не слишком плохой обусловленности.

Среднеквадратичные приближения и метод наименьших квадратов для решения переопределённых систем линейных алгебраических уравнений, которые возникают в связи с задачами обработки наблюдений, были почти одновременно предложены на рубеже XVIII–XIX веков А.М. Лежандром и К.Ф. Гауссом, причём первый дал новому подходу современное название. На практике метод наименьших квадратов очень часто применяется в силу двух главных причин. Во-первых, его применение бывает вызвано ясным содержательным смыслом задачи, в которой в качестве меры отклонения возникает именно сумма квадратов или интеграл от квадрата функции. К примеру, чрезвычайно популярно теоретико-вероятностное обоснование метода наименьших квадратов (см., к примеру, [49]). Впервые оно было дано также К.Ф. Гауссом и далее доведено до современного состояния в трудах П.С. Лапласа, П.Л. Чебышёва и потом А.А. Маркова и А.Н. Колмогорова. Во-вторых, в методе наименьших квадратов построение элемента наилучшего приближения сводится к решению системы линейных уравнений, т. е. хорошо разработанной задаче. Если для измерения расстояния между функциями применяются какие-то другие метрики, отличные от (2.100), то решение задачи минимизации этого расстояния может быть суще-

ственно более трудным. В целом, если какое-либо одно или оба из выписанных условий не выполняется, то метод наименьших квадратов может быть не самой лучшей возможностью решения задачи приближения.

Нередко форма приближающей функции (2.91) не подходит по тем или иным причинам, и тогда приходится прибегать к *нелинейному методу наименьших квадратов*, когда приближающая функция $g(x)$ выражается нелинейным образом через базисные функции $\varphi_1(x)$, $\varphi_2(x)$, \dots , $\varphi_m(x)$. Тогда минимизация среднеквадратичного отклонения f от g уже не сводится к решению системы линейных алгебраических уравнений (2.94), и для нахождения минимума нам нужно применять численные методы оптимизации. Обсуждение этого круга вопросов и дальнейшие ссылки можно найти, к примеру, в книге [41].

2.11 Полиномы Лежандра

2.11a Мотивация и определение

Примеры 2.10.2 и 2.10.3 из предшествующего раздела показывают, что выбор хорошего, т. е. ортогонального или почти ортогонального, базиса для среднеквадратичного приближения функций является нетривиальной задачей. Для её конструктивного решения можно воспользоваться, к примеру, известным из курса линейной алгебры процессом ортогонализации Грама-Шмидта или его модификациями (см. §3.7e). Напомним, что по данной конечной линейно независимой системе векторов v_1, v_2, \dots, v_n этот процесс строит ортогональный базис q_1, q_2, \dots, q_n линейной оболочки векторов v_1, v_2, \dots, v_n . Он имеет следующие расчётные формулы:

$$q_1 \leftarrow v_1, \quad (2.101)$$

$$q_k \leftarrow v_k - \sum_{i=1}^{k-1} \frac{\langle v_k, q_i \rangle}{\langle q_i, q_i \rangle} q_i, \quad k = 2, \dots, n. \quad (2.102)$$

Иногда получающийся ортогональный базис дополнительно нормируют.

В задаче среднеквадратичного приближения, рассмотренной в предшествующем §2.10g, ортогонализуемые элементы линейного пространства — это функции, а их скалярное произведение — интеграл (2.99). По

этой причине процесс ортогонализации (2.101)–(2.102) довольно трудоёмок, а конкретный вид функций, ортогональных в смысле $\mathcal{L}^2[a, b]$, которые получатся в результате, зависит, во-первых, от интервала $[a, b]$, для которого рассматривается скалярное произведение (2.99), и, во-вторых, от весовой функции $\varrho(x)$.

Для частного случая единичного веса, когда $\varrho(x) = 1$, мы можем существенно облегчить свою задачу, если найдём семейство ортогональных функций для какого-нибудь одного, канонического, интервала $[\alpha, \beta]$. Для любого другого интервала воспользуемся формулой линейной замены переменной $y = sx + r$ со специальным образом подобранными константами $r, s \in \mathbb{R}, s \neq 0$. Тогда $x = (y - r)/s$, и для $a = s\alpha + r, b = s\beta + r$ имеем равенство

$$\int_{\alpha}^{\beta} f(x) g(x) dx = \frac{1}{s} \int_a^b f\left(\frac{y-r}{s}\right) g\left(\frac{y-r}{s}\right) dy,$$

справедливое в силу формулы замены переменных в определённом интеграле. Из него вытекает, что равный нулю интеграл по каноническому интервалу $[\alpha, \beta]$ останется нулевым и при линейной замене переменных. Как следствие, получающиеся при такой замене функции $f\left(\frac{1}{s}(y-r)\right)$ и $g\left(\frac{1}{s}(y-r)\right)$ будут ортогональны на $[a, b]$.

Рассмотрим среднеквадратичное приближение функций полиномами. В этом случае в качестве канонического интервала $[\alpha, \beta]$ обычно берётся $[-1, 1]$, и тогда формула замены переменных принимает вид

$$y = \frac{1}{2}(b-a)x + \frac{1}{2}(a+b),$$

так что переменная y пробегает интервал $[a, b]$, если $x \in [-1, 1]$. Обратное преобразование даётся формулой

$$x = \frac{1}{b-a}(2y - (a+b)),$$

которая позволяет построить ортогональные в смысле $\mathcal{L}^2[a, b]$ полиномы для любого интервала $[a, b]$, зная их для $[-1, 1]$.

Полиномами Лежандра называют семейство полиномов $L_n(x)$, зависящих от неотрицательного целого параметра n , которые образуют ортогональную систему относительно скалярного произведения (2.99) с простейшим весом $\varrho(x) = 1$ на интервале $[-1, 1]$. Они были введены в широкий оборот французским математиком А. Лежандром в 1785 году. Из общей теории скалярного произведения в линейных пространствах

следует, что такие полиномы существуют и единственны с точностью до постоянного множителя. Нормирование полиномов Лежандра обычно выполняют различными способами, подходящими для той или иной задачи.

Применяя к степеням $1, x, x^2, x^3, \dots$ последовательно формулы ортогонализации (2.101)–(2.102) со скалярным произведением (2.99) на интервале $[-1, 1]$, получим

$$1, \quad x, \quad x^2 - \frac{1}{3}, \quad x^3 - \frac{3}{5}x, \quad \dots \quad (2.103)$$

(два первых полинома оказываются изначально ортогональными).

Часто в качестве альтернативного представления для полиномов Лежандра используют *формулу Родрига*

$$L_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n, \quad n = 0, 1, 2, \dots \quad (2.104)$$

Очевидно, что функция $L_n(x)$, определяемая этой формулой, является алгебраическим полиномом n -ой степени со старшим коэффициентом, не равным нулю, так как при n -кратном дифференцировании полинома $(x^2 - 1)^n = x^{2n} - nx^{2(n-1)} + \dots + (-1)^n$ степень понижается в точности на n . Коэффициент $1/(2^n n!)$ перед производной в (2.104) взят с той целью, чтобы удовлетворить условию $L_n(1) = 1$. Всюду далее посредством $L_n(x)$ мы будем обозначать полиномы Лежандра, определяемые формулой (2.104).

Предложение 2.11.1 *Полиномы $L_n(x)$, $n = 0, 1, \dots$, задаваемые формулой Родрига (2.104), ортогональны друг другу в смысле скалярного произведения на $\mathcal{L}^2[-1, 1]$ с единичным весом. Более точно,*

$$\int_{-1}^1 L_m(x)L_n(x) dx = \begin{cases} 0, & \text{если } m \neq n, \\ \frac{2}{2n+1}, & \text{если } m = n. \end{cases}$$

Доказательство. Обозначая

$$\psi(x) = (x^2 - 1)^n,$$

можно заметить, что

$$\psi^{(k)}(x) = \frac{d^k}{dx^k} (x^2 - 1)^n = 0 \quad \text{при } x = \pm 1, \quad k = 0, 1, 2, \dots, n-1.$$

Это следует из зануления множителей $(x^2 - 1)$, присутствующих во всех слагаемых выражений для $\psi^{(k)}(x)$, $k = 0, 1, \dots, n - 1$. Кроме того, в силу формулы Родрига (2.104)

$$L_n(x) = \frac{1}{2^n n!} \psi^{(n)}(x), \quad n = 0, 1, 2, \dots$$

Поэтому, если $Q(x)$ является n раз непрерывно дифференцируемой функцией на $[-1, 1]$, то, последовательно применяя n раз формулу интегрирования по частям, получим

$$\begin{aligned} \int_{-1}^1 Q(x) L_n(x) dx &= \frac{1}{2^n n!} \int_{-1}^1 Q(x) \psi^{(n)}(x) dx \\ &= \frac{1}{2^n n!} \left(Q(x) \psi^{(n-1)}(x) \right) \Big|_{-1}^1 - \frac{1}{2^n n!} \int_{-1}^1 Q'(x) \psi^{(n-1)}(x) dx \\ &= -\frac{1}{2^n n!} \int_{-1}^1 Q'(x) \psi^{(n-1)}(x) dx \\ &= \dots \\ &= (-1)^n \frac{1}{2^n n!} \int_{-1}^1 Q^{(n)}(x) \psi(x) dx. \end{aligned} \tag{2.105}$$

Если $Q(x)$ — любой полином степени меньше n , то его n -ая производная $Q^{(n)}(x)$ равна тождественному нулю, а потому из полученной формулы тогда следует

$$\int_{-1}^1 Q(x) L_n(x) dx = 0.$$

В частности, это верно и в случае, когда вместо $Q(x)$ берётся полином $L_m(x)$ степени m , меньшей n , что доказывает ортогональность этих полиномов с разными номерами.

Найдём теперь скалярное произведение полинома Лежандра с самим собой. Если $Q(x) = L_n(x)$, то

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

По этой причине из (2.105) следует

$$\begin{aligned} \int_{-1}^1 L_n(x) L_n(x) dx &= (-1)^n \frac{(2n)!}{(2^n n!)^2} \int_{-1}^1 \psi(x) dx \\ &= \frac{(2n)!}{(2^n n!)^2} \int_{-1}^1 (1-x^2)^n dx. \end{aligned}$$

Если обозначить

$$Z_n = \int_{-1}^1 (1-x^2)^n dx = \int_{-1}^1 (1-x^2)(1-x^2)^{n-1} dx,$$

то нетрудно найти, что

$$\begin{aligned} Z_n &= \int_{-1}^1 (1-x^2)^{n-1} dx - \int_{-1}^1 x^2(1-x^2)^{n-1} dx \\ &= Z_{n-1} - \int_{-1}^1 x^2(1-x^2)^{n-1} dx. \end{aligned}$$

Интегрируя по частям вычитаемое в последнем выражении, получим

$$\begin{aligned} \int_{-1}^1 x^2(1-x^2)^{n-1} dx &= -\frac{1}{2n} \int_{-1}^1 x d(1-x^2)^n \\ &= -\frac{1}{2n} \left(x(1-x^2)^n \right) \Big|_{-1}^1 + \frac{1}{2n} \int_{-1}^1 (1-x^2)^n dx = \frac{1}{2n} Z_n. \end{aligned}$$

Комбинируя эти результаты, будем иметь $Z_n = Z_{n-1} - \frac{1}{2n} Z_n$, откуда

$$Z_n = \frac{2n}{2n+1} Z_{n-1}.$$

Для нахождения числового значения Z_n учтём, что

$$Z_0 = \int_{-1}^1 (1-x^2)^0 dx = \int_{-1}^1 1 dx = 2.$$

Тогда $Z_1 = 2 \cdot \frac{2}{3}$ и т. д., так что

$$Z_n = 2 \frac{2 \cdot 4 \cdot \dots \cdot (2n-2) \cdot 2n}{3 \cdot 5 \cdot \dots \cdot (2n-1) \cdot (2n+1)}.$$

Окончательно, скалярное произведение $L_n(x)$ на себя равно

$$\frac{(2n)!}{(2^n n!)^2} Z_n = \frac{1 \cdot 2 \cdot 3 \cdot \dots \cdot 2n}{(2 \cdot 4 \cdot 6 \cdot \dots \cdot 2n)^2} Z_n = \frac{2}{2n+1}.$$

Это завершает доказательство предложения. ■

2.116 Основные свойства полиномов Лежандра

Выпишем первые полиномы Лежандра, как они даются формулой Родрига (2.104):

$$\begin{aligned} L_0(x) &= 1, \\ L_1(x) &= x, \\ L_2(x) &= \frac{1}{2}(3x^2 - 1), \\ L_3(x) &= \frac{1}{2}(5x^3 - 3x), \\ L_4(x) &= \frac{1}{8}(35x^4 - 30x^2 + 3), \\ L_5(x) &= \frac{1}{8}(63x^5 - 70x^3 + 15x), \\ &\dots \end{aligned} \tag{2.106}$$

Они с точностью до множителя совпадают с результатом ортогонализации Грама-Шмидта (2.103). Графики полиномов (2.106) изображены на Рис. 2.19, и они похожи на графики полиномов Чебышёва. В одном существенном моменте полиномы Лежандра всё же отличаются от полиномов Чебышёва: абсолютные значения локальных минимумов и максимумов на $[-1, 1]$ у полиномов Лежандра различны и не могут быть сделаны одинаковыми ни при каком масштабировании.

Тем не менее, сходство полиномов Лежандра и полиномов Чебышёва подтверждает следующее

Предложение 2.11.2 *Все нули полиномов Лежандра $L_n(x)$ вещественны, различны и находятся на интервале $[-1, 1]$.*

Доказательство. Предположим, что среди корней полинома $L_n(x)$, лежащих на $[-1, 1]$, имеется s штук различных корней $\theta_1, \theta_2, \dots, \theta_s$

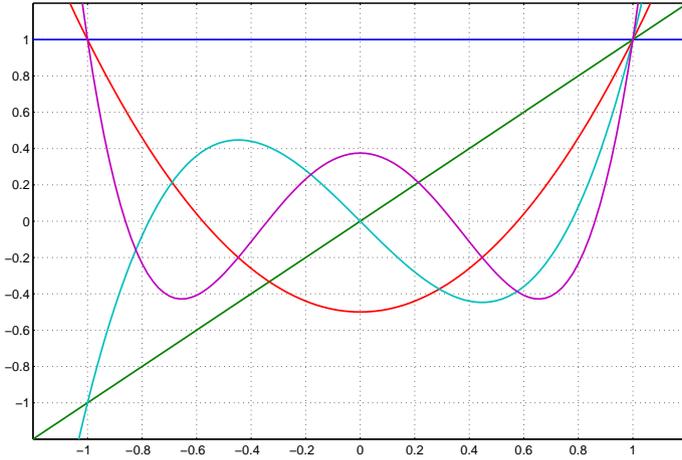


Рис. 2.19. Графики первых полиномов Лежандра на интервале $[-1.2, 1.2]$.

нечётной кратности $\alpha_1, \alpha_2, \dots, \alpha_s$, так что

$$L_n(x) = (x - \theta_1)^{\alpha_1} (x - \theta_2)^{\alpha_2} \cdots (x - \theta_s)^{\alpha_s} \gamma(x),$$

где в полиноме $\gamma(x)$ присутствуют корни $L_n(x)$, не лежащие на $[-1, 1]$, а также те корни $L_n(x)$ из $[-1, 1]$, которые имеют чётную кратность. Таким образом, $\gamma(x)$ уже не меняет знака на интервале $[-1, 1]$. Ясно, что $s \leq n$, и наша задача — установить равенство $s = n$.

Рассмотрим интеграл

$$\begin{aligned} \mathcal{I} &= \int_{-1}^1 L_n(x) (x - \theta_1)(x - \theta_2) \cdots (x - \theta_s) dx \\ &= \int_{-1}^1 (x - \theta_1)^{\alpha_1+1} (x - \theta_2)^{\alpha_2+1} \cdots (x - \theta_s)^{\alpha_s+1} \gamma(x) dx. \end{aligned}$$

Теперь $\alpha_1+1, \alpha_2+1, \dots, \alpha_s+1$ — чётные числа, так что подинтегральное выражение не меняет знак на $[-1, 1]$. Это выражение равно нулю лишь в конечном множестве точек, и потому определён $\mathcal{I} \neq 0$.

С другой стороны, выражение для \mathcal{I} есть скалярное произведение, в смысле $\mathcal{L}^2[-1, 1]$, полинома $L_n(x)$ на полином $(x - \theta_1)(x - \theta_2) \cdots (x - \theta_s)$

степени не более $n - 1$, если выполнено условие $s < n$. Следовательно, в силу свойств полиномов Лежандра при этом должно быть $\mathcal{I} = 0$.

Полученное противоречие может быть снято только в случае $s = n$, т. е. когда $\mathcal{I} \neq 0$. При этом все корни полинома $L_n(x)$ различны и лежат на интервале $[-1, 1]$. ■

Отметим, что проведённое доказательство проходит для скалярных произведений вида (2.99) с достаточно произвольными весовыми функциями $\varrho(x)$, а не только для единичного веса. Кроме того, тот факт, что интервал интегрирования есть $[-1, 1]$, также нигде не использовался в явном виде. Фактически, это доказательство годится даже для бесконечных пределов интегрирования. Оно показывает, что корни любых ортогональных полиномов вещественны и различны.

Можно показать дополнительно, что нули полинома Лежандра $L_n(x)$ перемежаются с нулями полинома $L_{n+1}(x)$. Наконец, аналогично полиномам Чебышёва, нули полиномов Лежандра также сгущаются к концам интервала $[-1, 1]$.

Ещё одно интересное свойство полиномов Лежандра, задаваемых посредством формулы Родрига (2.104):

$$L_n(1) = 1, \quad L_n(-1) = (-1)^n, \quad n = 0, 1, 2, \dots$$

Кроме того, справедливо рекуррентное представление

$$(n + 1)L_{n+1}(x) = (2n + 1)xL_n(x) - nL_{n-1}(x).$$

Доказательство этих свойств можно найти, в частности, в [25, 53]. Последняя формула даёт практически удобный способ вычисления значений полиномов Лежандра, так как в их явном представлении (2.106) коэффициенты растут экспоненциально быстро в зависимости от номера полинома и, как следствие, прямые вычисления с ними могут дать большую погрешность.

Введём так называемые *приведённые полиномы Лежандра* $\tilde{L}_n(x)$, старший коэффициент у которых равен единице. Как следствие формулы Родрига (2.104) можем выписать следующее представление:

$$\tilde{L}_n(x) = \frac{1}{2n(2n-1)\cdots(n+1)} \frac{d^n}{dx^n} (x^2 - 1)^n = \frac{n!}{(2n)!} \frac{d^n}{dx^n} (x^2 - 1)^n,$$

$n = 1, 2, \dots$ Как и исходная формула Родрига, выражение после второго равенства имеет также смысл при $n = 0$, если под производной нулевого порядка от функции понимать её саму.

Предложение 2.11.3 Среди всех полиномов степени n , $n \geq 1$, со старшим коэффициентом, равным 1, полином $\tilde{L}_n(x)$ имеет на интервале $[-1, 1]$ наименьшее среднеквадратичное отклонение от нуля. Иными словами, если $Q_n(x)$ — полином степени n со старшим коэффициентом 1, то

$$\int_{-1}^1 (Q_n(x))^2 dx \geq \int_{-1}^1 (\tilde{L}_n(x))^2 dx. \quad (2.107)$$

Доказательство. Если $Q_n(x) = x^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0$, то для отыскания наименьшего значения выражения

$$\begin{aligned} \mathcal{J}(a_0, a_1, \dots, a_{n-1}) &= \int_{-1}^1 (Q_n(x))^2 dx \\ &= \int_{-1}^1 (x^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0)^2 dx \end{aligned} \quad (2.108)$$

продифференцируем этот интеграл по коэффициентам a_0, a_1, \dots, a_{n-1} и приравняем полученные производные к нулю. Так как в данных условиях дифференцирование интеграла по параметру, от которого зависит подинтегральная функция, сводится к взятию интеграла от производной, то имеем в результате

$$\begin{aligned} \frac{\partial \mathcal{J}}{\partial a_k} &= \int_{-1}^1 2(x^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0) x^k dx \\ &= 2 \int_{-1}^1 Q_n(x) x^k dx = 0, \end{aligned}$$

$k = 0, 1, \dots, n-1$. Это означает, что полином $Q_n(x)$ ортогонален в смысле $\mathcal{L}^2[-1, 1]$ всем полиномам меньшей степени. Следовательно, при минимальном значении интеграла (2.108) полином $Q_n(x)$ обязан совпадать с n -ым полиномом Лежандра. ■

Для построения полинома, который имеет наименьшее среднеквадратичное отклонение от нуля на произвольном интервале $[a, b]$ можно воспользоваться линейной заменой переменной и затем масштабированием, аналогично тому, как это было сделано для полиномов Чебышёва в §2.3б.

Помимо полиномов Лежандра существуют и другие семейства ортогональных полиномов, широко используемые в теории и практических вычислениях. В частности, введённые в §2.3 полиномы Чебышёва образуют семейство полиномов, ортогональных на интервале $[-1, 1]$ с весом $(1 - x^2)^{-1/2}$.

Часто возникает необходимость воспользоваться ортогональными полиномами на бесконечных интервалах $[0, +\infty]$ или даже $[-\infty, \infty]$. Естественно, единичный вес $\varrho(x) = 1$ тут малоприменим, так как с ним интегралы по бесконечным интервалам окажутся, по большей части, расходящимися. Полиномы, ортогональные на интервалах $[0, +\infty]$ или $[-\infty, \infty]$ с быстроубывающими весами e^{-x} и e^{-x^2} называются полиномами Лагерра и полиномами Эрмита соответственно.¹⁴ Они также находят многообразные применения в задачах приближения, и более подробные сведения на эту тему читатель может почерпнуть в [25, 63].

2.12 Численное интегрирование

2.12a Постановка и обсуждение задачи

Задача вычисления определённого интеграла

$$\int_a^b f(x) dx \quad (2.109)$$

является одной из важнейших математических задач, к которой сводится большое количество различных вопросов теории и практики. Это нахождение площадей криволинейных фигур, центров тяжести и моментов инерции тел, работы переменной силы и т. п. механические, физические, химические и другие задачи. В математическом анализе обобщается *формула Ньютона-Лейбница*

$$\int_a^b f(x) dx = F(b) - F(a), \quad (2.110)$$

где $F(x)$ — первообразная для функции $f(x)$, т. е. такая, что $F'(x) = f(x)$. Она даёт удобный способ вычисления интегралов, который в значительной степени удовлетворяет потребности решения подобных задач. Тем не менее, возникают ситуации, когда для вычисления интеграла (2.109) требуются другие подходы.

¹⁴Иногда их называют также полиномами Чебышёва-Лагерра и Чебышёва-Эрмита (см., к примеру, [39, 63]), поскольку они были известны ещё П.Л. Чебышёву.

Задачей *численного интегрирования* называют задачу нахождения определённого интеграла (2.109) на основе знания значений функции $f(x)$, без привлечения её первообразных и формулы Ньютона-Лейбница (2.110). Подобная задача нередко возникает на практике, например, если подынтегральная функция $f(x)$ задана таблично, т. е. своими значениями в дискретном наборе точек, а не аналитической формулой. В некоторых случаях численное нахождение интеграла приходится выполнять потому, что первообразная для интегрируемой функции не выражается через элементарные функции. Но даже если эта первообразная может быть найдена в конечном виде, её вычисление не всегда осуществляется просто (длинное и неустойчивое к ошибкам округления выражение и т. п.). Все эти причины вызывают необходимость развития численных методов для нахождения определённых интегралов.

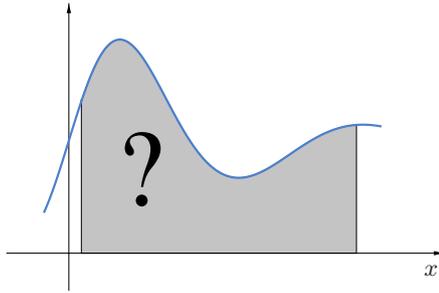


Рис. 2.20. Вычисление определённого интеграла необходимо при нахождении площадей фигур с криволинейными границами

Наибольшее распространение в вычислительной практике получили формулы вида

$$\int_a^b f(x) dx \approx \sum_{k=0}^n c_k f(x_k), \quad (2.111)$$

где c_k — некоторые постоянные коэффициенты, x_k — точки из интервала интегрирования $[a, b]$, $k = 0, 1, \dots, n$. Подобные формулы называют *квадратурными формулами*,¹⁵ коэффициенты c_k — это *весовые коэффициенты* или просто *веса* квадратурной формулы, а точки x_k — её *узлы*.

¹⁵ «Квадратура» в оригинальном смысле, восходящем ещё к античности, означала построение квадрата, равновеликого заданной фигуре. Но в эпоху Возрождения этот термин стал означать вычисление площадей фигур.

В многомерном случае аналогичные приближённые равенства

$$\int_D f(x) dx \approx \sum_{k=0}^n c_k f(x_k),$$

где $x, x_k \in D \subset \mathbb{R}^m$, D — область в \mathbb{R}^m , $m \geq 2$,

называют *кубатурными формулами*. Естественное условие принадлежности узлов x_k области интегрирования вызвано тем, что за её пределами подинтегральная функция может быть просто не определена, как, например, $\arcsin x$ вне интервала $[-1, 1]$.

Тот факт, что квадратурные и кубатурные формулы являются линейными выражениями от значений интегрируемой функции в узлах, объясняется линейным характером зависимости самого интеграла от подинтегральной функции. С другой стороны, квадратурные формулы можно рассматривать как обобщения интегральных сумм Римана (через которые интеграл Римана и определяется), так как простейшие составные квадратурные формулы прямоугольников просто совпадают с этими интегральными суммами.

Как и ранее, совокупность узлов x_0, x_1, \dots, x_n квадратурной (кубатурной) формулы называют *сеткой*. Разность

$$R(f) = \int_a^b f(x) dx - \sum_{k=0}^n c_k f(x_k)$$

называется *погрешностью квадратурной формулы* или её *остаточным членом*. Это число, зависящее от подинтегральной функции f , в отличие от остаточного члена интерполяции, который является ещё функцией точки.

Если для некоторой функции f или же для целого класса функций $\mathcal{F} \ni f$ имеет место точное равенство

$$\int_a^b f(x) dx = \sum_{k=0}^n c_k f(x_k),$$

то будем говорить, что квадратурная формула *точна* (является точной) на f или для целого класса функций \mathcal{F} . То, насколько широким является класс функций, на котором точна рассматриваемая формула, может служить косвенным признаком её точности вообще. Очень часто в качестве класса «пробных функций» \mathcal{F} , для которых исследуется

совпадение результата квадратурной формулы и искомого интеграла, берут алгебраические полиномы. В этой связи полезно

Определение 2.12.1 *Алгебраической степенью точности квадратурной формулы называют наибольшую степень алгебраических полиномов, для которых эта квадратурная формула является точной.*

Соответственно, из двух квадратурных формул более предпочтительной будем считать ту, которая имеет бóльшую алгебраическую степень точности. Неформальным обоснованием этого критерия служит тот факт, что с помощью полиномов более высокой степени можно получать более точные приближения функций, как локально (с помощью формулы Тейлора), так и глобально (к примеру, с помощью разложения по полиномам Чебышёва или Лежандра).

Рассмотрим теперь влияние погрешностей реальных вычислений на ответ, получаемый с помощью квадратурных формул. Предположим, что значения $f(x_k)$ интегрируемой функции в узлах x_k вычисляются неточно, с погрешностями δ_k . Тогда при вычислениях по квадратурной формуле получим

$$\sum_{k=0}^n c_k (f(x_k) + \delta_k) = \sum_{k=0}^n c_k f(x_k) + \sum_{k=0}^n c_k \delta_k.$$

Если для всех $k = 0, 1, \dots, n$ знаки погрешностей δ_k совпадают со знаками весов c_k , то общая абсолютная погрешность результата, полученного по квадратурной формуле, становится равной $\sum_k |c_k| \delta_k$. Следовательно, сумму модулей весов квадратурной формулы, т. е. величину

$$\sum_{k=0}^n |c_k|,$$

нужно рассматривать как коэффициент усиления погрешности при вычислениях с этой формулой.

Если при значительном количестве узлов n мы хотим организовать вычисления по квадратурной формуле наиболее устойчивым образом, то все весовые коэффициенты c_k должны быть положительны: именно тогда при прочих равных условиях сумма модулей весов минимальна. В частности, в случае интегрирования функций, принимающих значения одного знака, мы избегаем тогда потери точности при вычитании близких значений, которое могло бы случиться в формуле, где одновременно присутствуют положительные и отрицательные веса.

2.126 Формулы Ньютона-Котеса. Простейшие квадратурные формулы

Простейший приём построения квадратурных формул — замена подинтегральной функции $f(x)$ на интервале интегрирования $[a, b]$ на «более простую», легче интегрируемую функцию, которая интерполирует или приближает $f(x)$ по заданным узлам x_0, x_1, \dots, x_n . В случае, когда $f(x)$ заменяется интерполянтom и все рассматриваемые узлы — простые, говорят о квадратурных формулах интерполяционного типа, или, что равносильно, об *интерполяционных квадратурных формулах*. Наиболее часто подинтегральную функцию интерполируют полиномами, и в нашем курсе мы будем рассматривать только такие интерполяционные квадратурные формулы.

Формулами Ньютона-Котеса называют интерполяционные квадратурные формулы, полученные с помощью алгебраической интерполяции подинтегральной функции на равномерной сетке с простыми узлами. В зависимости от того, включаются ли концы интервала интегрирования $[a, b]$ в множество узлов квадратурной формулы или нет, различают формулы Ньютона-Котеса *замкнутого типа* и *открытого типа*.

Далее мы построим и исследуем формулы Ньютона-Котеса для $n = 0, 1, 2$, причём будем строить наиболее популярные формулы замкнутого типа за исключением случая $n = 0$, когда имеем всего один узел и замкнутая формула просто невозможна.

Если $n = 0$, то подинтегральная функция $f(x)$ интерполируется полиномом нулевой степени, т. е. какой-то константой, равной значению $f(x)$ в единственном узле x_0 . Если взять $x_0 = a$, то получается квадратурная формула «левых прямоугольников», а если $x_0 = b$ — формула «правых прямоугольников» (см. Рис. 2.21).

Ещё один естественный вариант выбора единственного узла —

$$x_0 = \frac{1}{2}(a + b),$$

т. е. как середины интервала интегрирования $[a, b]$. При этом приходим к квадратурной формуле

$$\int_a^b f(x) dx \approx (b - a) \cdot f\left(\frac{a + b}{2}\right),$$

называемой *формулой средних прямоугольников*: согласно ей интеграл берётся равным площади прямоугольника с основанием $(b - a)$ и высо-

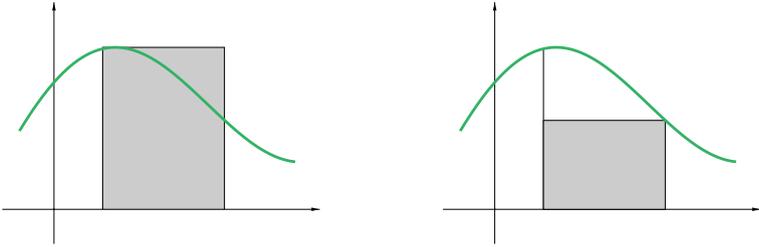


Рис. 2.21. Иллюстрация квадратурных формул левых и правых прямоугольников

той $f((a+b)/2)$ (см. Рис. 2.22). Эту формулу нередко называют также просто «формулой прямоугольников», так как она является наиболее часто используемым вариантом рассмотренных простейших квадратурных формул.

Оценим погрешность формулы средних прямоугольников методом локальных разложений, который ранее был использован при исследовании численного дифференцирования. Разлагая $f(x)$ в окрестности точки $x_0 = \frac{1}{2}(a+b)$ по формуле Тейлора с точностью до членов первого порядка, получим

$$f(x) = f\left(\frac{a+b}{2}\right) + f'\left(\frac{a+b}{2}\right) \cdot \left(x - \frac{a+b}{2}\right) + \frac{f''(\xi)}{2} \cdot \left(x - \frac{a+b}{2}\right)^2,$$

где ξ — зависящая от x точка интервала $[a, b]$, которую корректно обозначить через $\xi(x)$. Далее

$$\begin{aligned} R(f) &= \int_a^b f(x) dx - (b-a) \cdot f\left(\frac{a+b}{2}\right) \\ &= \int_a^b \left(f(x) - f\left(\frac{a+b}{2}\right) \right) dx \\ &= \int_a^b \left(f'\left(\frac{a+b}{2}\right) \cdot \left(x - \frac{a+b}{2}\right) + \frac{f''(\xi(x))}{2} \cdot \left(x - \frac{a+b}{2}\right)^2 \right) dx \\ &= \int_a^b \frac{f''(\xi(x))}{2} \cdot \left(x - \frac{a+b}{2}\right)^2 dx, \end{aligned}$$

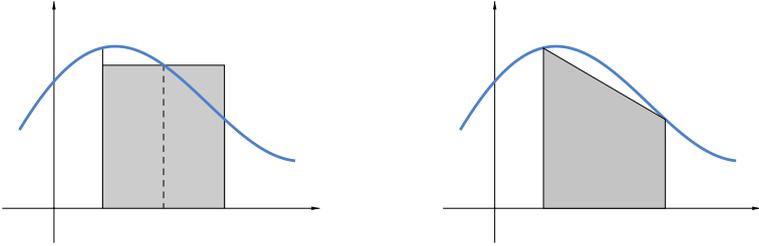


Рис. 2.22. Иллюстрация квадратурных формул средних прямоугольников и трапеций

поскольку

$$\int_a^b \left(x - \frac{a+b}{2}\right) dx = \int_{-\frac{b-a}{2}}^{\frac{b-a}{2}} t dt = 0,$$

— интеграл от первого члена разложения зануляется. Следовательно, с учётом принятого нами ранее обозначения

$$M_p := \max_{x \in [a,b]} |f^{(p)}(x)|$$

можно выписать оценку

$$\begin{aligned} |R(f)| &\leq \int_a^b \left| \frac{f''(\xi)}{2} \right| \cdot \left(x - \frac{a+b}{2}\right)^2 dx \leq \frac{M_2}{2} \int_a^b \left(x - \frac{a+b}{2}\right)^2 dx \\ &= \frac{M_2}{2} \cdot \frac{1}{3} \left(x - \frac{a+b}{2}\right)^3 \Big|_a^b = \frac{M_2(b-a)^3}{24}. \end{aligned}$$

Отсюда, в частности, следует, что для полиномов степени не выше 1 формула (средних) прямоугольников даёт точное значение интеграла, коль скоро вторая производная подинтегральной функции тогда зануляется и $M_2 = 0$.

Полученная оценка точности неулучшаема, так как достигается на функции $g(x) = \left(x - \frac{1}{2}(a+b)\right)^2$. При этом

$$M_2 = \max_{x \in [a,b]} |g''(x)| = 2, \quad g\left(\frac{a+b}{2}\right) = 0,$$

и потому

$$\int_a^b g(x) dx - (b-a) \cdot g\left(\frac{a+b}{2}\right) = \frac{(b-a)^3}{12} = \frac{M_2(b-a)^3}{24},$$

т. е. имеем точное равенство на погрешность.

Рассмотрим теперь квадратурную формулу Ньютона-Котеса, соответствующую случаю $n = 1$, когда подынтегральная функция приближается интерполяционным полиномом первой степени. Построим его по узлам $x_0 = a$ и $x_1 = b$:

$$P_1(x) = \frac{x-b}{a-b} f(a) + \frac{x-a}{b-a} f(b).$$

Интегрируя это равенство, получим

$$\begin{aligned} \int_a^b P_1(x) dx &= \frac{f(a)}{a-b} \int_a^b (x-b) dx + \frac{f(b)}{b-a} \int_a^b (x-a) dx \\ &= \frac{f(a)}{a-b} \frac{(x-b)^2}{2} \Big|_a^b + \frac{f(b)}{b-a} \frac{(x-a)^2}{2} \Big|_a^b \\ &= \frac{b-a}{2} (f(a) + f(b)). \end{aligned}$$

Мы вывели *квадратурную формулу трапеций*

$$\int_a^b f(x) dx \approx (b-a) \cdot \frac{f(a) + f(b)}{2}, \quad (2.112)$$

название которой также навеяно геометрическим образом. Фактически, согласно этой формуле точное значение интеграла заменяется на значение площади трапеции (стоящей боком на оси абсцисс) с высотой $(b-a)$ и основаниями, равными $f(a)$ и $f(b)$ (см. Рис. 2.22).

Чтобы найти погрешность формулы трапеций, вспомним оценку (2.24) для погрешности интерполяционного полинома. Из неё следует, что

$$f(x) - P_1(x) = \frac{f''(\xi(x))}{2} \cdot (x-a)(x-b)$$

для некоторой точки $\xi(x) \in [a, b]$. Таким образом, для формулы трапеций

$$R(f) = \int_a^b (f(x) - P_1(x)) dx = \int_a^b \frac{f''(\xi(x))}{2} \cdot (x-a)(x-b) dx,$$

но вычисление полученного интеграла на практике нереально из-за неизвестного вида $\xi(x)$. Как обычно, имеет смысл вывести какие-то более удобные оценки погрешности, хотя они, возможно, будут не столь точны.

Поскольку выражение $(x-a)(x-b)$ почти всюду на интервале $[a, b]$ сохраняет один и тот же знак, то

$$\begin{aligned} |R(f)| &\leq \int_a^b \frac{|f''(\xi(x))|}{2} \cdot |(x-a)(x-b)| dx \\ &\leq \frac{M_2}{2} \cdot \left| \int_a^b (x-a)(x-b) dx \right|, \end{aligned}$$

где $M_2 = \max_{x \in [a, b]} |f''(x)|$. Далее

$$\begin{aligned} \int_a^b (x-a)(x-b) dx &= \int_a^b (x^2 - (a+b)x + ab) dx \\ &= \frac{x^3}{3} \Big|_a^b - (a+b) \frac{x^2}{2} \Big|_a^b + abx \Big|_a^b \quad (2.113) \\ &= \frac{1}{6} \left(2(b^3 - a^3) - 3(a+b)(b^2 - a^2) + 6ab(b-a) \right) \\ &= \frac{1}{6} (-b^3 + 3ab^2 - 3a^2b + a^3) = -\frac{(b-a)^3}{6}. \end{aligned}$$

Поэтому окончательно

$$|R(f)| \leq \frac{M_2(b-a)^3}{12}.$$

Эта оценка погрешности квадратурной формулы трапеций неулучшаема, поскольку достигается на функции $g(x) = (x-a)^2$.

2.12в Квадратурная формула Симпсона

Построим квадратурную формулу Ньютона-Котеса для $n = 2$, т. е. для трёх равномерно расположенных узлов

$$x_0 = a, \quad x_1 = \frac{a+b}{2}, \quad x_2 = b$$

из интервала интегрирования $[a, b]$.

Для упрощения рассуждений выполним параллельный перенос криволинейной трапеции, площадь которой мы находим с помощью интегрирования, и сделаем точку a началом координат оси абсцисс (см. Рис. 2.23). Тогда правым концом интервала интегрирования станет $l = b - a$. Пусть

$$\check{P}_2(x) = c_0 + c_1x + c_2x^2$$

— полином второй степени, интерполирующий сдвинутую подинтегральную функцию по узлам 0 , $l/2$ и l . Если график $\check{P}_2(x)$ проходит через точки плоскости с координатами

$$(0, f(a)), \quad \left(\frac{l}{2}, f\left(\frac{a+b}{2}\right) \right), \quad (l, f(b)),$$

то

$$\begin{cases} c_0 = f(a), \\ c_0 + c_1 \frac{l}{2} + c_2 \frac{l^2}{4} = f\left(\frac{a+b}{2}\right), \\ c_0 + c_1 l + c_2 l^2 = f(b). \end{cases} \quad (2.114)$$

Площадь, ограниченная графиком интерполяционного полинома, равна

$$\begin{aligned} \int_0^l (c_0 + c_1x + c_2x^2) dx &= c_0l + c_1 \frac{l^2}{2} + c_2 \frac{l^3}{3} \\ &= \frac{l}{6} (6c_0 + 3c_1l + 2c_2l^2). \end{aligned}$$

Фактически, для построения квадратурной формулы требуется решить относительно c_0 , c_1 и c_2 систему уравнений (2.114) и потом подставить результаты в полученное выше выражение. Но можно выразить трёхчлен $6c_0 + 3c_1l + 2c_2l^2$ через значения подинтегральной функции f в узлах, не решая систему (2.114) явно.

Умножая второе уравнение системы (2.114) на 4 и складывая с первым и третьим уравнением, получим

$$6c_0 + 3c_1l + 2c_2l^2 = f(a) + 4f\left(\frac{a+b}{2}\right) + f(b).$$

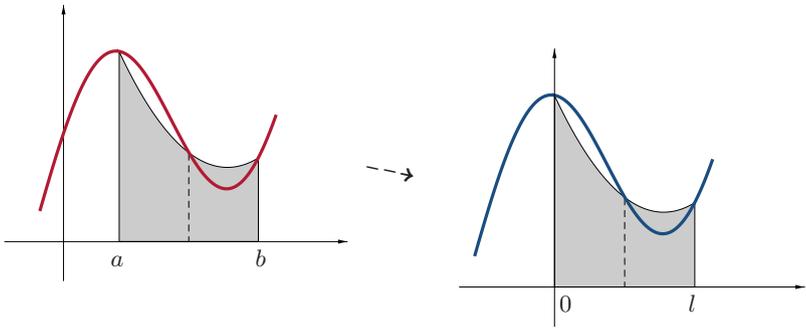


Рис. 2.23. Иллюстрация вывода квадратурной формулы Симпсона

Таким образом,

$$\begin{aligned} \int_0^l \check{P}_2(x) dx &= \int_0^l (c_0 + c_1x + c_2x^2) dx \\ &= \frac{l}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right), \end{aligned}$$

что даёт приближённое равенство

$$\int_a^b f(x) dx \approx \frac{b-a}{6} \cdot \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right). \quad (2.115)$$

Оно называется *квадратурной формулой Симпсона* или *формулой парабол* (см. Рис. 2.23), коль скоро основано на приближении подинтегральной функции подходящей параболой.¹⁶

Читатель может самостоятельно убедиться, что та же самая формула получается в результате интегрирования по $[a, b]$ интерполяционного

¹⁶Приведённый нами элегантный вывод формулы Симпсона восходит к учебнику А.Н. Крылова [15].

полинома второй степени в форме Лагранжа

$$\begin{aligned}
 P_2(x) &= \frac{\left(x - \frac{a+b}{2}\right)(x-b)}{\left(a - \frac{a+b}{2}\right)(a-b)} f(a) + \frac{(x-a)(x-b)}{\left(\frac{a+b}{2} - a\right)\left(\frac{a+b}{2} - b\right)} f\left(\frac{a+b}{2}\right) \\
 &\quad + \frac{(x-a)\left(x - \frac{a+b}{2}\right)}{(b-a)\left(b - \frac{a+b}{2}\right)} f(b) \\
 &= \frac{2}{(b-a)^2} \left(\left(x - \frac{a+b}{2}\right)(x-b) f(a) - 2(x-a)(x-b) f\left(\frac{a+b}{2}\right) \right. \\
 &\quad \left. + (x-a)\left(x - \frac{a+b}{2}\right) f(b) \right),
 \end{aligned}$$

который строится для подынтегральной функции по узлам a , $(a+b)/2$ и b .

Предложение 2.12.1 *Квадратурная формула Симпсона имеет алгебраическую степень точности 3, т. е. является точной для любого полинома степени не выше третьей.*

Доказательство. Отметим прежде всего, что для полиномов степени не выше второй этот факт следует прямо из того, что формула Симпсона построена как интерполяционная квадратурная формула, основанная на интерполяции подынтегральной функции полиномом второй степени. Поэтому достаточно показать, что формула Симпсона точна для монома x^3 , но не является точной для более высоких степеней.

Имеем

$$\int_a^b x^3 dx = \frac{b^4 - a^4}{4}.$$

С другой стороны,

$$\begin{aligned} \frac{b-a}{6} \left(a^3 + 4 \left(\frac{a+b}{2} \right)^3 + b^3 \right) &= \frac{b-a}{6} \left(a^3 + \frac{a^3 + 3a^2b + 3ab^2 + b^3}{2} + b^3 \right) \\ &= \frac{b-a}{6} \cdot \frac{3a^3 + 3a^2b + 3ab^2 + 3b^3}{2} \\ &= \frac{b-a}{4} (a^3 + a^2b + ab^2 + b^3) = \frac{b^4 - a^4}{4}, \end{aligned}$$

что совпадает с результатом точного интегрирования.

Для монома x^4 длинными, но несложными выкладками нетрудно проверить, что результат, даваемый формулой Симпсона для интеграла по интервалу $[a, b]$, т. е.

$$\frac{b-a}{6} \left(a^4 + 4 \left(\frac{a+b}{2} \right)^4 + b^4 \right),$$

отличается от точного значения интеграла

$$\int_a^b x^4 dx = \frac{b^5 - a^5}{5}$$

на величину $(b-a)^5/120$. Она не зануляется при $a \neq b$, так что на полиномах четвёртой степени формула Симпсона уже не точна. ■

Итак, несмотря на то, что формула Симпсона основана на интерполяции подинтегральной функции полиномом степени 2, фактическая точность формулы более высока, чем та, что обеспечивается полиномом второй степени. В этой ситуации для более аккуратной оценки погрешности формулы Симпсона на основе известной погрешности алгебраической интерполяции (аналогично выводу погрешности формулы трапеций в §2.12б) желательно взять более высокую степень переменной в выражении для погрешности. Иными словами, при оценке погрешности формулы Симпсона нужно взять для подинтегральной функции интерполяционный полином третьей степени. При наличии всего трёх узлов мы находимся в условиях задачи интерполяции с кратными узлами.

Предполагая существование производной f' в среднем узле $x_1 = (a+b)/2$, можно считать, к примеру, что именно он является кратным

узлом. Формально мы будем решать задачу построения такого интерполяционного полинома 3-й степени $H_3(x)$, что

$$H_3(a) = f(a), \quad H_3(b) = f(b), \quad (2.116)$$

$$H_3\left(\frac{a+b}{2}\right) = f\left(\frac{a+b}{2}\right), \quad H_3'\left(\frac{a+b}{2}\right) = f'\left(\frac{a+b}{2}\right), \quad (2.117)$$

Хотя конкретное значение производной в средней точке $(a+b)/2$ далее никак не будет использоваться. Здесь нам важно лишь то, что при любом значении этой производной решение задачи (2.116)–(2.117) существует, и потребуется оценка его погрешности.

Существование и единственность решения подобных задач была установлена в §2.4 и там же обосновывается оценка его погрешности (2.44):

$$f(x) - H_m(x) = \frac{f^{(m+1)}(\xi(x))}{(m+1)!} \prod_{i=0}^n (x - x_i)^{N_i}, \quad (2.118)$$

где N_i — кратности узлов, $m = N_0 + N_1 + \dots + N_n - 1$ — степень интерполяционного полинома, а $\xi(x)$ — некоторая точка из $[a, b]$, зависящая от x . Для решения задачи (2.116)–(2.117) справедливо

$$f(x) - H_3(x) = \frac{f^{(4)}(\xi(x))}{24} \cdot (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b).$$

Далее, из того, что формула Симпсона точна для полиномов третьей степени, а также из условий (2.116)–(2.117) следуют равенства

$$\begin{aligned} \int_a^b H_3(x) dx &= \frac{b-a}{6} \cdot \left(H_3(a) + 4H_3\left(\frac{a+b}{2}\right) + H_3(b) \right) \\ &= \frac{b-a}{6} \cdot \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right). \end{aligned} \quad (2.119)$$

Отсюда уже нетрудно вывести выражение для погрешности квадра-

турной формулы Симпсона:

$$\begin{aligned}
 R(f) &= \int_a^b f(x) dx - \frac{b-a}{6} \cdot \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right) \\
 &= \int_a^b (f(x) - H_3(x)) dx \quad \text{в силу (2.119)} \\
 &= \int_a^b \frac{f^{(4)}(\xi(x))}{24} \cdot (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b) dx \quad \text{из (2.118)}.
 \end{aligned}$$

Из него следует оценка

$$\begin{aligned}
 |R(f)| &= \left| \int_a^b \frac{f^{(4)}(\xi(x))}{24} \cdot (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b) dx \right| \\
 &\leq \int_a^b \left| \frac{f^{(4)}(\xi(x))}{24} \right| \cdot \left| (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b) \right| dx \\
 &\leq \frac{M_4}{24} \cdot \left| \int_a^b (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b) dx \right|, \quad (2.120)
 \end{aligned}$$

коль скоро в интегрируемой функции подвыражение $(x-a)(x - (a+b)/2)^2(x-b)$ не меняет знак на интервале интегрирования $[a, b]$. Здесь обозначено $M_4 = \max_{x \in [a, b]} |f^{(4)}(x)|$.

Для вычисления фигурирующего в (2.120) интеграла сделаем замену переменных

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

тогда

$$\begin{aligned} & \int_a^b (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b) dx \\ &= \int_{-\frac{b-a}{2}}^{\frac{b-a}{2}} \left(t + \frac{b-a}{2}\right) t^2 \left(t - \frac{b-a}{2}\right) dt \\ &= \int_{-\frac{b-a}{2}}^{\frac{b-a}{2}} t^2 \left(t^2 - \frac{(b-a)^2}{4}\right) dt = -\frac{(b-a)^5}{120}. \end{aligned}$$

Окончательно

$$|R(f)| \leq \frac{M_4 (b-a)^5}{2880}.$$

Как видим, более тонкие рассуждения о свойствах формулы Симпсона позволили получить действительно более точную оценку её погрешности.

2.12г Общие интерполяционные квадратурные формулы

Квадратурными формулами интерполяционного типа мы назвали (см. §2.12б) формулы, получающиеся в результате замены подынтегральной функции $f(x)$ интерполяционным полиномом $P_n(x)$, который построен по некоторой совокупности простых узлов x_0, x_1, \dots, x_n из интервала интегрирования. Выпишем для общего случая этот полином в форме Лагранжа:

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

где

$$\phi_i(x) = \frac{(x-x_0) \cdots (x-x_{i-1})(x-x_{i+1}) \cdots (x-x_{n+1})}{(x_i-x_0) \cdots (x_i-x_{i-1})(x_i-x_{i+1}) \cdots (x_i-x_{n+1})},$$

— базисные полиномы Лагранжа (стр. 51).

Интерполяционная квадратурная формула должна получаться из приближённого равенства

$$\int_a^b f(x) dx \approx \int_a^b P_n(x) dx \quad (2.121)$$

в результате выполнения интегрирования в правой части. Как следствие, в представлении (2.111) весовые коэффициенты формулы имеют вид

$$\begin{aligned} c_i &= \int_a^b \phi_i(x) dx \\ &= \int_a^b \frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_{n+1})}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_{n+1})} dx, \end{aligned} \quad (2.122)$$

$i = 0, 1, \dots, n$. Эти значения весов c_i , определяемых по узлам x_0, x_1, \dots, x_n , являются отличительным характеристическим признаком именно интерполяционной квадратурной формулы. Если для заданного набора узлов у какой-либо квадратурной формулы весовые коэффициенты равны (2.122), то можно считать, что она построена на основе алгебраической интерполяции подинтегральной функции по этим узлам, взятым с единичной кратностью.

Теорема 2.12.1 *Для того, чтобы квадратурная формула (2.111), построенная по $(n + 1)$ попарно различным узлам, была интерполяционной, необходимо и достаточно, чтобы её алгебраическая степень точности была не меньшей n .*

В качестве замечания к формулировке нужно отметить, что в условиях теоремы квадратурная формула на самом деле может иметь алгебраическую степень точности выше n , как, например, формула средних прямоугольников или формула Симпсона.

Доказательство. Необходимость условий теоремы очевидна: интерполяционная квадратурная формула на $n + 1$ узлах, конечно же, точна на полиномах степени n , поскольку тогда подинтегральная функция совпадает со своим алгебраическим интерполянтном.

Покажем достаточность: если квадратурная формула (2.111), построенная по $(n + 1)$ узлу, является точной для любого алгебраического полинома степени n , то её весовые коэффициенты вычисляются по

формулам (2.122), т. е. она является квадратурной формулой интерполяционного типа.

В самом деле, для базисных интерполяционных полиномов $\phi_i(x)$ выполнено свойство (2.9)

$$\phi_i(x_j) = \delta_{ij} = \begin{cases} 0, & \text{при } i \neq j, \\ 1, & \text{при } i = j, \end{cases}$$

и они имеют степень n . Следовательно, применяя рассматриваемую квадратурную формулу для вычисления интеграла от $\phi_i(x)$, получим

$$\int_a^b \phi_i(x) dx = \sum_{k=0}^n c_k \phi_i(x_k) = \sum_{k=0}^n c_k \delta_{ik} = c_i,$$

и это верно для всех $i = 0, 1, \dots, n$. Иными словами, имеет место равенство (2.122), что и требовалось доказать. ■

В частности, если в интерполяционной квадратурной формуле вместо подынтегральной функции взять полином $P_0(x) = x^0 = 1$, то получаем равенство

$$b - a = \int_a^b 1 dx = \sum_{k=0}^n c_k,$$

— сумма весов такой квадратурной формулы равна длине интервала интегрирования.

Из (2.121) ясно, что погрешность интерполяционных квадратурных формул равна

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

где $R_n(f, x)$ — остаточный член алгебраической интерполяции. В §2.2 была получена оценка для $R_n(f, x)$ в форме Коши (2.24)

$$R_n(f, x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \cdot \omega_n(x),$$

где $\xi(x) \in [a, b]$, и поэтому

$$R(f) = \frac{1}{(n+1)!} \int_a^b f^{(n+1)}(\xi(x)) \omega_n(x) dx.$$

Справедлива огрублённая оценка

$$|R(f)| \leq \frac{M_{n+1}}{(n+1)!} \int_a^b |\omega_n(x)| dx, \quad (2.123)$$

где $M_{n+1} = \max_{x \in [a,b]} |f^{(n+1)}(x)|$. Из неё можно ещё раз заключить, что квадратурная формула интерполяционного типа, построенная по $(n+1)$ узлам, является точной для любого полинома степени не более n , поскольку тогда $M_{n+1} = 0$.

Из наших рассуждений видно, что оценка (2.123) является простейшей, использующей лишь основные свойства алгебраического интерполанта. В некоторых случаях она может оказаться существенно завышенной, как это имеет место, к примеру, для формулы Симпсона.

2.12д Дальнейшие формулы Ньютона-Котеса

В §2.12б и §2.12в простейшие квадратурные формулы Ньютона-Котеса — формулы прямоугольников и трапеций, формула Симпсона — были выведены и исследованы средствами, индивидуальными для каждой отдельной формулы. В этом разделе мы взглянем на формулы Ньютона-Котеса с более общих позиций.

Зафиксировав номер n , $n \geq 1$, возьмём на интервале интегрирования $[a, b]$ равноотстоящие друг от друга узлы

$$x_k^{(n)} = a + kh, \quad k = 0, 1, \dots, n, \quad h = \frac{b-a}{n}.$$

Для определения весов формул Ньютона-Котеса необходимо вычислить величины (2.122), которые мы обозначим для рассматриваемого частного случая как

$$A_k^{(n)} = \int_a^b \frac{(x - x_0^{(n)}) \cdots (x - x_{k-1}^{(n)}) (x - x_{k+1}^{(n)}) \cdots (x - x_{n+1}^{(n)})}{(x_k^{(n)} - x_0^{(n)}) \cdots (x_k^{(n)} - x_{k-1}^{(n)}) (x_k^{(n)} - x_{k+1}^{(n)}) \cdots (x_k^{(n)} - x_n^{(n)})} dx,$$

$k = 0, 1, \dots, n$. Сделаем в этом интеграле замену переменных $x = a + th$,

где t пробегает интервал $[0, n]$. Тогда

$$\begin{aligned} dx &= h dt, \\ (x - x_0^{(n)}) \cdots (x - x_{k-1}^{(n)}) (x - x_{k+1}^{(n)}) \cdots (x - x_{n+1}^{(n)}) \\ &= h^n t(t-1) \cdots (t-k+1)(t-k-1) \cdots (t-n), \\ (x_k^{(n)} - x_0^{(n)}) \cdots (x_k^{(n)} - x_{k-1}^{(n)}) (x_k^{(n)} - x_{k+1}^{(n)}) \cdots (x_k^{(n)} - x_n^{(n)}) \\ &= (-1)^{n-k} h^n k!(n-k)!, \end{aligned}$$

где считается, что $0! = 1$. Окончательно

$$A_k^{(n)} = h \frac{(-1)^{n-k}}{k!(n-k)!} \int_0^n t(t-1) \cdots (t-k+1)(t-k-1) \cdots (t-n) dt,$$

$k = 0, 1, \dots, n$. Чтобы придать результату не зависящий от интервала интегрирования вид, положим

$$A_k^{(n)} = (b-a) B_k^{(n)},$$

где

$$B_k^{(n)} = \frac{(-1)^{n-k}}{n k!(n-k)!} \int_0^n t(t-1) \cdots (t-k+1)(t-k-1) \cdots (t-n) dt.$$

Теперь уже величины $B_k^{(n)}$ не зависят от h и $[a, b]$. Они называются *коэффициентами Котеса*.

К примеру, для $n = 1$

$$B_0^{(1)} = - \int_0^1 (t-1) dt = - \left. \frac{(t-1)^2}{2} \right|_0^1 = \frac{1}{2},$$

$$B_1^{(1)} = \int_0^1 t dt = \left. \frac{t^2}{2} \right|_0^1 = \frac{1}{2}.$$

Мы вновь получили веса квадратурной формулы трапеций (2.112). Для

случая $n = 2$

$$B_0^{(2)} = \frac{1}{4} \int_0^2 (t-1)(t-2) dt = \frac{1}{4} \left(\frac{t^3}{3} - 3\frac{t^2}{2} + 2t \right) \Big|_0^2 = \frac{1}{6},$$

$$B_1^{(2)} = -\frac{1}{2} \int_0^2 t(t-2) dt = -\frac{1}{2} \left(\frac{t^3}{3} - t^2 \right) \Big|_0^2 = \frac{4}{6},$$

$$B_2^{(2)} = \frac{1}{4} \int_0^2 t(t-1) dt = \frac{1}{4} \left(\frac{t^3}{3} - \frac{t^2}{2} \right) \Big|_0^2 = \frac{1}{6}.$$

Полученные коэффициенты соответствуют формуле Симпсона (2.115). И так далее.

За прошедшие три столетия коэффициенты Котеса были тщательно вычислены для значений n из начального отрезка натурального ряда. В Табл. 2.2, заимствованной из книги [15], приведены коэффициенты Котеса для $n \leq 10$ (см. также [3, 21, 35, 67]).

Можно видеть, что с ростом n значения коэффициентов Котеса $B_k^{(n)}$ в зависимости от номера k начинают всё сильнее и сильнее «осциллировать» (напоминая в чём-то пример Рунге, стр. 84). Результатом этого является то необычное и противоестественное обстоятельство, что среди весов формул Ньютона-Котеса при числе узлов $n = 8, 10$ и больших встречаются отрицательные. Это снижает ценность соответствующих формул, так как при интегрировании знакопостоянных функций может приводить к вычитанию близких чисел и потере точности.

К середине XX века выяснилось, что отмеченный недостаток типичен для формул Ньютона-Котеса высоких порядков. Р.О. Кузьмин получил в [48] асимптотические формулы для коэффициентов Котеса¹⁷, из которых следует, что сумма их модулей, т. е.

$$\sum_{k=0}^n |B_k^{(n)}|,$$

неограниченно возрастает с ростом n . Отсюда вытекает, во-первых, что погрешности вычислений с формулами Ньютона-Котеса могут быть

¹⁷Помимо оригинальной статьи Р.О. Кузьмина [48] эти асимптотические формулы можно также увидеть в учебнике [17].

Таблица 2.2. Коэффициенты Котеса

n	1	2	3	4	5	6	7	8	9	10
$k=0$	1	$\frac{1}{6}$	$\frac{1}{8}$	$\frac{7}{90}$	$\frac{19}{288}$	$\frac{41}{840}$	$\frac{751}{17280}$	$\frac{989}{28350}$	$\frac{2857}{89600}$	$\frac{16067}{598752}$
$k=1$	1	$\frac{4}{6}$	$\frac{3}{8}$	$\frac{16}{45}$	$\frac{25}{96}$	$\frac{9}{35}$	$\frac{3577}{17280}$	$\frac{5838}{28350}$	$\frac{15741}{89600}$	$\frac{106300}{598752}$
$k=2$		$\frac{1}{6}$	$\frac{3}{8}$	$\frac{2}{15}$	$\frac{25}{144}$	$\frac{9}{280}$	$\frac{1323}{17280}$	$-\frac{928}{28350}$	$\frac{1080}{89600}$	$-\frac{48525}{598752}$
$k=3$			$\frac{1}{8}$	$\frac{16}{45}$	$\frac{25}{144}$	$\frac{34}{105}$	$\frac{2989}{17280}$	$\frac{10496}{28350}$	$\frac{19344}{89600}$	$\frac{272400}{598752}$
$k=4$				$\frac{7}{90}$	$\frac{25}{96}$	$\frac{9}{280}$	$\frac{2989}{17280}$	$-\frac{4540}{28350}$	$\frac{5778}{89600}$	$-\frac{260550}{598752}$
$k=5$					$\frac{19}{288}$	$\frac{9}{35}$	$\frac{1323}{17280}$	$\frac{10496}{28350}$	$\frac{5778}{89600}$	$\frac{427368}{598752}$
$k=6$						$\frac{41}{840}$	$\frac{3577}{17280}$	$-\frac{928}{28350}$	$\frac{19344}{89600}$	$-\frac{260550}{598752}$
$k=7$							$\frac{751}{17280}$	$\frac{5838}{28350}$	$\frac{1080}{89600}$	$\frac{272400}{598752}$
$k=8$								$\frac{989}{28350}$	$\frac{15741}{89600}$	$-\frac{48525}{598752}$
$k=9$									$\frac{2857}{89600}$	$\frac{106300}{598752}$
$k=10$										$\frac{16067}{598752}$

сколь угодно велики (см. §2.12а). Во-вторых, так как дополнительно

$$\sum_{k=0}^n B_k^{(n)} = \frac{1}{b-a} \sum_{k=0}^n A_k^{(n)} = \frac{1}{b-a} \int_a^b 1 \, dx = 1,$$

то при достаточно больших n среди коэффициентов $B_k^{(n)}$ обязательно должны быть как положительные, так и отрицательные. Доказательство упрощённого варианта этого результата можно найти в [25].

Общую теорию квадратурных формул Ньютона-Котеса вместе с тщательным исследованием их погрешностей читатель может увидеть, к примеру, в книгах [3, 17, 53]. Следует сказать, что формулы Ньютона-

Котеса высоких порядков не очень употребительны. Помимо отмеченной выше численной неустойчивости они проигрывают по точности результатов на одинаковом количестве узлов формулам Гаусса (изучаемым далее в §2.13) и другим квадратурным формулам.

Из популярных квадратурных формул Ньютона-Котеса приведём ещё формулу «трёх восьмых», которая получается при замене подинтегральной функции интерполяционным полиномом 3-й степени:

$$\int_a^b f(x) dx \approx \frac{b-a}{8} \cdot \left(f(a) + 3f\left(\frac{2a+b}{3}\right) + 3f\left(\frac{a+2b}{3}\right) + f(b) \right). \quad (2.124)$$

Её погрешность оценивается как

$$|R(f)| \leq \frac{M_4(b-a)^5}{6480},$$

где $M_4 = \max_{x \in [a,b]} |f^{(4)}(x)|$. Порядок точности этой формулы — такой же, как и у формулы Симпсона. Вообще, можно показать, что формулы Ньютона-Котеса с нечётным числом узлов, один из которых приходится на середину интервала интегрирования, имеют (как формула Симпсона) повышенный порядок точности [1, 3, 17].

2.12e Метод неопределённых коэффициентов

Пусть требуется вычислить интеграл

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

для которого мы ищем приближение в виде

$$\int_a^b f(x) dx \approx \sum_{k=0}^n c_k f(x_k),$$

с заданными узлами x_0, x_1, \dots, x_n . Весовые коэффициенты c_k можно найти из условия зануления погрешности этого равенства для какого-то «достаточно представительного» набора несложно интегрируемых функций $f_i(x)$, $i = 1, 2, \dots$. Каждое отдельное равенство является уравнением на неизвестные c_0, c_1, \dots, c_n , и потому, выписав достаточное

число подобных уравнений и решив полученную систему, мы сможем определить желаемые веса, т. е. построить квадратурную формулу. В этом — суть *метода неопределённых коэффициентов*. Он идейно похож, таким образом, на метод неопределённых коэффициентов для построения формул численного дифференцирования из §2.8в.

В качестве пробных функций $f_p(x)$, $p = 1, 2, \dots$ часто берут алгебраические полиномы. Для равномерно расположенных узлов при этом получаются знакомые нам квадратурные формулы Ньютона-Котеса.

Продемонстрируем работу метода неопределённых коэффициентов для тригонометрических полиномов

$$1, \quad \sin(px), \quad \cos(px), \quad p = 1, 2, \dots$$

2.13 Квадратурные формулы Гаусса

2.13а Задача оптимизации квадратур

Параметрами квадратурной формулы (2.111)

$$\int_a^b f(x) dx \approx \sum_{k=0}^n c_k f(x_k)$$

являются узлы x_k и весовые коэффициенты c_k , $k = 0, 1, \dots, n$. Однако, строя квадратурные формулы Ньютона-Котеса, мы заранее задавали положение узлов, равномерное на интервале интегрирования, и потом по ним находили веса. Таким образом, возможности общей формулы (2.111) были использованы не в полной мере, поскольку для достижения наилучших результатов можно было бы управлять ещё и положением узлов. Лишь в формуле средних прямоугольников положение единственного узла было выбрано из соображений симметрии, и это привело к существенному улучшению точности. Напомним для примера, что специальное неравномерное расположение узлов интерполяции по корням полиномов Чебышёва существенно улучшает точность интегрирования (см. §2.3).

Здесь, правда, возникает весьма важный методический вопрос: как измерять это «улучшение» квадратурной формулы? Что брать критерием её точности? В идеальном случае желательно было бы минимизировать погрешность квадратурной формулы для тех или иных классов функций, но в такой общей постановке задача делается довольно

сложной (хотя и не неразрешимой). Один из возможных естественных ответов на поставленный вопрос состоит в том, чтобы в качестве меры того, насколько хороша и точна квадратурная формула, брать её алгебраическую степень точности (см. Определение 2.12.1).

Как следствие, сформулированную в начале этого параграфа задачу оптимизации узлов можно поставить, к примеру, следующим образом: для заданного фиксированного числа узлов из интервала интегрирования нужно построить квадратурную формулу, которая имеет наивысшую алгебраическую степень точности, т.е. является точной на полиномах наиболее высокой степени. Нетривиальное решение этой задачи действительно существует, и формулы наивысшей алгебраической степени точности называются *квадратурными формулами Гаусса*, поскольку впервые они были рассмотрены в начале XIX века К.Ф. Гауссом.

Далее для удобства мы будем записывать квадратурные формулы Гаусса не в виде (2.111), а как

$$\int_a^b f(x) dx \approx \sum_{k=1}^n c_k f(x_k), \quad (2.125)$$

нумеруя узлы с $k = 1$, а не с нуля. Требование точного равенства для любого полинома степени m в этой формуле в силу её линейности эквивалентно тому, что формула является точной для одночленов $f(x) = x^l$, $l = 0, 1, 2, \dots, m$, т.е.

$$\int_a^b x^l dx = \sum_{k=1}^n c_k x_k^l, \quad l = 0, 1, 2, \dots, m,$$

или

$$\sum_{k=1}^n c_k x_k^l = \frac{1}{l+1} (b^{l+1} - a^{l+1}), \quad l = 0, 1, 2, \dots, m. \quad (2.126)$$

Это система из $(m+1)$ нелинейных уравнений с $2n$ неизвестными величинами $c_1, c_2, \dots, c_n, x_1, x_2, \dots, x_n$. Число уравнений совпадает с числом неизвестных при $m+1 = 2n$, т.е. $m = 2n - 1$, и это, вообще говоря, есть максимальное возможное значение для m . При больших значениях m система уравнений (2.126) переопределена и в случае общего положения оказывается неразрешимой.

Сделанное заключение можно обосновать строго.

Предложение 2.13.1 *Алгебраическая степень точности квадратурной формулы, построенной по n узлам, не может превосходить $2n-1$.*

Доказательство. Пусть x_1, x_2, \dots, x_n — узлы квадратурной формулы (2.125). Рассмотрим интегрирование по интервалу $[a, b]$ функции

$$g(x) = ((x - x_1)(x - x_2) \cdots (x - x_n))^2,$$

которая является полиномом степени $2n$. Если квадратурная формула (2.125) точна для $g(x)$, то

$$\sum_{k=1}^n c_k g(x_k) = \sum_{k=1}^n c_k \cdot 0 = 0,$$

тогда как значение интеграла от $g(x)$ очевидно не равно нулю. Подынтегральная функция $g(x)$ всюду на $[a, b]$ положительна за исключением лишь конечного множества точек — узлов x_1, x_2, \dots, x_n , и поэтому $\int_a^b g(x) dx > 0$.

Полученное противоречие показывает, что квадратурная формула (2.125) не является точной для полиномов степени $2n$. ■

Итак, наивысшая алгебраическая степень точности квадратурной формулы, построенной по n узлам, в общем случае может быть равна $2n-1$. Для двух узлов это 3, при трёх узлах имеем 5, и т. д. Для сравнения напомним, что алгебраические степени точности формул трапеций и Симпсона, построенных по двум и трём узлам соответственно, равны всего 1 и 3. При возрастании числа узлов этот выигрыш в алгебраической степени точности формул Гаусса, достигаемый за счёт разумного расположения узлов, нарастает.

2.136 Простейшие квадратуры Гаусса

Перейдём к построению квадратурных формул Гаусса. При небольших n система уравнений (2.126) для узлов и весов может быть решена с помощью несложных аналитических преобразований.

Пусть $n = 1$, тогда $m = 2n - 1 = 1$, и система уравнений (2.126) принимает вид

$$\begin{cases} c_1 = b - a, \\ c_1 x_1 = \frac{1}{2}(b^2 - a^2). \end{cases}$$

Отсюда

$$c_1 = b - a,$$

$$x_1 = \frac{1}{2c_1}(b^2 - a^2) = \frac{1}{2}(a + b).$$

Как легко видеть, получающаяся квадратурная формула — это формула (средних) прямоугольников

$$\int_a^b f(x) dx \approx (b - a) \cdot f\left(\frac{a + b}{2}\right).$$

Нам в самом деле известно (см. §2.126), что она резко выделяется своей точностью среди родственных квадратурных формул.

Пусть $n = 2$, тогда алгебраическая степень точности соответствующей квадратурной формулы равна $m = 2n - 1 = 3$. Система уравнений (2.126) для узлов и весов принимает вид

$$\begin{cases} c_1 + c_2 = b - a, \\ c_1 x_1 + c_2 x_2 = \frac{1}{2}(b^2 - a^2), \\ c_1 x_1^2 + c_2 x_2^2 = \frac{1}{3}(b^3 - a^3), \\ c_1 x_1^3 + c_2 x_2^3 = \frac{1}{4}(b^4 - a^4). \end{cases}$$

Она обладает определённой симметрией: одновременная перемена местами x_1 с x_2 и c_1 с c_2 оставляет систему неизменной. По этой причине, учитывая вид первого уравнения, будем искать решение, в котором $c_1 = c_2$. Это даёт

$$c_1 = c_2 = \frac{1}{2}(b - a),$$

и из первого и второго уравнений тогда получаем

$$x_1 + x_2 = a + b. \quad (2.127)$$

Отсюда после возведения в квадрат имеем

$$x_1^2 + 2x_1x_2 + x_2^2 = a^2 + 2ab + b^2. \quad (2.128)$$

В то же время, с учётом найденных значений c_1 и c_2 из третьего уравнения системы следует

$$x_1^2 + x_2^2 = \frac{2}{3}(b^2 + ab + a^2),$$

и, вычитая это равенство из (2.128), будем иметь

$$x_1 x_2 = \frac{1}{6}(b^2 + 4ab + a^2). \quad (2.129)$$

Соотношения (2.127) и (2.129) на основе известной из элементарной алгебры теоремы Виета позволяют сделать вывод, что x_1 и x_2 являются корнями квадратного уравнения

$$x^2 - (a + b)x + \frac{1}{6}(b^2 + 4ab + a^2) = 0,$$

так что

$$x_{1,2} = \frac{1}{2}(a + b) \pm \frac{\sqrt{3}}{6}(b - a). \quad (2.130)$$

Удовлетворение полученными решениями четвёртого уравнения системы проверяется прямой подстановкой. Кроме того, поскольку

$$\frac{1}{2} > \frac{1}{2} \cdot \frac{1}{\sqrt{3}} = \frac{\sqrt{3}}{6},$$

то x_1 и x_2 действительно лежат на интервале $[a, b]$. В целом мы вывели квадратурную формулу Гаусса

$$\int_a^b f(x) dx = \frac{b - a}{2} \cdot (f(x_1) + f(x_2)), \quad (2.131)$$

где узлы x_1 и x_2 определяются посредством (2.130).

Пример 2.13.1 Вычислим с помощью полученной выше формулы Гаусса с двумя узлами (2.131) интеграл

$$\int_0^{\pi/2} \cos x dx,$$

точное значение которого согласно формуле Ньютона-Лейбница равно $\sin(\pi/2) - \sin 0 = 1$. В соответствии с (2.130) и (2.131) имеем

$$\begin{aligned} \int_0^{\pi/2} \cos x dx &\approx \frac{\pi/2}{2} \cdot \left(\cos\left(\frac{\pi}{4} - \frac{\sqrt{3}}{6} \frac{\pi}{2}\right) + \cos\left(\frac{\pi}{4} + \frac{\sqrt{3}}{6} \frac{\pi}{2}\right) \right) \\ &= 0.998473. \end{aligned}$$

Формула Ньютона-Котеса с двумя узлами 0 и $\pi/2$ — формула трапеций — даёт для этого интеграла значение

$$\int_0^{\pi/2} \cos x \, dx \approx \frac{\pi}{2} \cdot \left(\cos 0 + \cos \frac{\pi}{2} \right) = 0.785398,$$

точность которого весьма низка.

Чтобы получить с формулами Ньютона-Котеса точность вычисления рассматриваемого интеграла, сравнимую с той, что даёт формула Гаусса, придется брать больше узлов. Так, формула Симпсона (2.115), использующая три узла — 0, $\pi/4$ и $\pi/2$, — приводит к результату

$$\begin{aligned} \int_0^{\pi/2} \cos x \, dx &\approx \frac{\pi/2}{6} \cdot \left(\cos 0 + 4 \cos \frac{\pi}{4} + \cos \frac{\pi}{2} \right) \\ &= \frac{\pi}{12} (1 + 2\sqrt{2}) = 1.00228, \end{aligned}$$

погрешность которого по порядку величины примерно равна погрешности ответа по формуле Гаусса (2.131), но всё-таки превосходит её в полтора раза. ■

С ростом n сложность системы уравнений (2.126) для узлов и весов формул Гаусса быстро нарастает, так что в общем случае остаётся неясным, будет ли разрешима эта система (2.126) при любом наперёд заданном n . Будут ли её решения вещественными? Сколько их всего? Будут ли эти решения принадлежать интервалу $[a, b]$, чтобы служить практически удобными узлами квадратурной формулы?

Получение ответов на поставленные вопросы непосредственно из системы уравнений (2.126) представляется громоздким и малоперспективным. К.Ф. Гауссом было предложено расчленив получившуюся задачу на отдельные подзадачи

- 1) построения узлов формулы и
- 2) вычисления её весовых коэффициентов.

Зная узлы формулы, можно подставить их в систему уравнений (2.126), которая в результате решительно упростится, превратившись в систему линейных алгебраических уравнений относительно c_1, c_2, \dots, c_n . Она будет переопределённой, но нам достаточно рассматривать подсистему из первых n уравнений, матрица которой является матрицей

Вандермонда относительно узлов x_1, x_2, \dots, x_n . Решение этой подсистемы даст искомые веса квадратурной формулы Гаусса. Можно показать, что они будут также удовлетворять оставшимся n уравнениям системы (2.126) (см., к примеру, [9]).

Другой способ решения подзадачи 2 — вычисление весовых коэффициентов по формулам (2.122), путём интегрирования коэффициентов интерполяционного полинома Лагранжа. В этом случае мы пользуемся тем фактом, что конструируемая квадратурная формула Гаусса оказывается квадратурной формулой интерполяционного типа. Это прямо следует из Теоремы 2.12.1, коль скоро формула Гаусса, построенная по n узлам, является точной для полиномов степени не менее $n - 1$. Детали этого построения и конкретные выкладки читатель может найти, к примеру, в [3].

2.13в Выбор узлов для квадратурных формул Гаусса

Теорема 2.13.1 *Квадратурная формула (2.125)*

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

имеет алгебраическую степень точности $(2n - 1)$ тогда и только тогда, когда

- (1) *она является интерполяционной квадратурной формулой;*
- (2) *её узлы x_1, x_2, \dots, x_n суть корни такого полинома*

$$\omega(x) = (x - x_1)(x - x_2) \cdots (x - x_n),$$

что

$$\int_a^b \omega(x) q(x) dx = 0 \tag{2.132}$$

для любого полинома $q(x)$ степени не выше $(n - 1)$.

Выражение

$$\int_a^b \omega(x) q(x) dx$$

— интеграл от произведения двух функций, уже встречалось нам в §2.10г. Мы могли видеть, что на пространстве $\mathcal{L}^2[a, b]$ всех интегрируемых с квадратом функций оно задаёт *скалярное произведение*, т. е. симметричную билинейную и положительно определённую форму. По этой причине утверждение Теоремы 2.13.1 часто формулируют так: для того, чтобы квадратурная формула

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

построенная по n узлам, имела алгебраическую степень точности $(2n - 1)$, необходимо и достаточно, чтобы эта формула была интерполяционной, а её узлы x_1, x_2, \dots, x_n являлись корнями полинома $\omega(x)$, который с единичным весом ортогонален на $[a, b]$ любому полиному степени не выше $(n - 1)$.

Доказательство. Необходимость. Пусть рассматриваемая квадратурная формула имеет алгебраическую степень точности $(2n - 1)$, т. е. точна на полиномах степени $(2n - 1)$. Таковым является, в частности, полином $\omega(x)q(x)$, имеющий степень не выше $n + (n - 1)$, если степень $q(x)$ не превосходит $(n - 1)$. Тогда справедливо точное равенство

$$\int_a^b \omega(x)q(x) dx = \sum_{k=1}^n c_k \omega(x_k)q(x_k) = 0,$$

поскольку все $\omega(x_k) = 0$. Так как этот результат верен для любого полинома $q(x)$ степени не выше $n - 1$, то отсюда следует выполнение условия (2).

Справедливость условия (1) следует из Теоремы 2.12.1: если построенная по n узлам квадратурная формула (2.125) является точной для любого полинома степени не менее $n - 1$, то она — интерполяционная.

Достаточность. Пусть имеется полином $\omega(x)$ степени n , имеющий n различных корней на интервале $[a, b]$ и удовлетворяющий условию ортогональности (2.132) с любым полиномом $q(x)$ степени не выше $(n - 1)$. Покажем, что интерполяционная квадратурная формула, построенная по узлам x_1, x_2, \dots, x_n , которые являются корнями $\omega(x)$, будет точна на полиномах степени $2n - 1$.

Пусть $f(x)$ — произвольный полином степени $2n - 1$. Тогда после деления его на $\omega(x)$ получим представление

$$f(x) = \omega(x)q(x) + r(x), \tag{2.133}$$

где $q(x)$ и $r(x)$ — соответственно частное и остаток от деления $f(x)$ на $\omega(x)$. При этом полином $q(x)$ имеет степень $(2n - 1) - n = n - 1$, а степень полинома-остатка $r(x)$ по определению меньше степени $\omega(x)$, т. е. не превосходит $n - 1$. Отсюда

$$\int_a^b f(x) dx = \int_a^b \omega(x) q(x) dx + \int_a^b r(x) dx = \int_a^b r(x) dx \quad (2.134)$$

в силу сделанного нами предположения об ортогональности $\omega(x)$ всем полиномам степени не выше $n - 1$.

Но по условиям теоремы рассматриваемая квадратурная формула является интерполяционной и построена по n узлам. Поэтому она является точной на полиномах степени $n - 1$ (см. Теорему 2.12.1), в частности, на полиноме $r(x)$. Следовательно,

$$\begin{aligned} \int_a^b r(x) dx &= \sum_{k=1}^n c_k r(x_k) = \sum_{k=1}^n c_k (\omega(x_k) q(x_k) + r(x_k)) \\ &\quad \text{в силу равенств } \omega(x_k) = 0 \\ &= \sum_{k=1}^n c_k f(x_k), \quad \text{поскольку имеет место (2.133)}. \end{aligned}$$

Итак, сравнивая результаты этой выкладки с (2.134), будем иметь

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

т. е. исследуемая квадратурная формула действительно является точной на полиномах степени $2n - 1$. ■

Подведём промежуточные итоги. Процедура построения квадратурных формул Гаусса разделена нами на две отдельные задачи нахождения узлов и вычисления весов. В свою очередь, узлы квадратурной формулы, как выясняется, можно взять корнями некоторых специальных полиномов $\omega(x)$, удовлетворяющих условиям Теоремы 2.13.1. В этих полиномах легко угадываются знакомые нам из §2.11 ортогональные полиномы, которые являются полиномами Лежандра для случая $[a, b] = [-1, 1]$ или соответствующим образом преобразованы из них для произвольного интервала интегрирования $[a, b]$.

2.13г Практическое применение формул Гаусса

Отдельное нахождение узлов и весов формул Гаусса для каждого конкретного интервала интегрирования $[a, b]$ является весьма трудозатратным, и если бы нам нужно было проделывать эту процедуру всякий раз при смене интервала $[a, b]$, то практическое применение формул Гаусса значительно потеряло бы свою привлекательность. Естественная идея состоит в том, чтобы найти узлы и веса формул Гаусса для какого-то одного «канонического» интервала, а затем получать их для любого другого интервала с помощью несложных преобразований.

В качестве канонического интервала обычно берут $[-1, 1]$, т. е. тот интервал, для которого строятся ортогональные полиномы Лежандра. Этот интервал также удобен симметричностью относительно нуля, которая позволяет более просто использовать свойство симметрии узлов и весовых коэффициентов квадратурной формулы. В §2.11 мы указали рецепт построения из полиномов Лежандра полиномов, ортогональных с единичным весом, для любого интервала вещественной оси. Этой техникой и нужно воспользоваться в данном случае.

Если

$$x = \frac{1}{2}(a + b) + \frac{1}{2}(b - a)y, \quad (2.135)$$

то переменная x будет пробегать интервал $[a, b]$, когда y изменяется в $[-1, 1]$. Обратное преобразование даётся формулой

$$y = \frac{1}{b - a}(2x - (a + b)).$$

В частности, если $y_i, i = 1, 2, \dots, n$, — корни полинома Лежандра, которые согласно Предложению 2.11.2 все различны и лежат на интервале $[-1, 1]$, то узлы квадратурной формулы Гаусса для интервала интегрирования $[a, b]$ суть

$$x_i = \frac{1}{2}(a + b) + \frac{1}{2}(b - a)y_i, \quad i = 1, 2, \dots, n. \quad (2.136)$$

Все они также различны и лежат на интервале интегрирования $[a, b]$.

Далее, веса c_k любой интерполяционной квадратурной формулы могут быть выражены в виде интегралов (2.122). В случае формул Гаусса (когда узлы нумеруются с единицы) они принимают вид

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

где $\phi_k(x)$ — k -ый базисный полином Лагранжа (см. стр. 51), построенный по узлам (2.136):

$$\phi_k(x) = \frac{(x - x_1) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_1) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)}.$$

Тогда, выполняя замену переменных (2.135), получим

$$dx = d\left(\frac{1}{2}(a+b) + \frac{1}{2}(b-a)y\right) = \frac{1}{2}(b-a) dy,$$

и потому

$$c_k = \int_a^b \phi_k(x) dx = \frac{1}{2}(b-a) \int_{-1}^1 \phi_k(y) dy, \quad k = 1, 2, \dots, n,$$

где $\phi_k(y)$ — k -ый базисный полином Лагранжа, построенный по узлам $y_i, i = 1, 2, \dots, n$, которые являются корнями n -го полинома Лежандра. Получается, что веса квадратурной формулы Гаусса для произвольного интервала интегрирования $[a, b]$ вычисляются простым умножением весов для канонического интервала $[-1, 1]$ на множитель $\frac{1}{2}(b-a)$ — радиус интервала интегрирования.

Для интервала $[-1, 1]$ узлы квадратурных формул Гаусса (т. е. корни полиномов Лежандра) и их веса тщательно заатабулированы для первых натуральных чисел n вплоть до нескольких десятков. Обсуждение вычислительных формул и других деталей численных процедур для их нахождения читатель может найти, к примеру, в книгах [3, 53] и в специальных журнальных статьях. В частности, оказывается, что весовые коэффициенты формулы Гаусса с n узлами даются выражением

$$c_k = \frac{2}{(1 - x_k^2)(L'_n(x_k))^2}, \quad k = 1, 2, \dots, n,$$

где $L_n(x)$ — n -ый полином Лежандра в форме, даваемой формулой Родрига (2.104).

Конкретные числовые значения узлов и весов квадратур Гаусса приводятся в подробных руководствах по вычислительным методам [2, 3, 9, 15, 16, 56] или в специализированных справочниках, например, в [35, 47]. В частности, в учебнике [3] значения весов и узлов формул Гаусса приведены для небольших n с 16 значащими цифрами, в книге [16] — с 15 значащими цифрами вплоть до $n = 16$, а в справочниках

Таблица 2.3. Узлы и веса квадратурных формул Гаусса

Узлы	Веса
$n = 2$	
$\pm 0.57735\ 02691\ 89626$	$1.00000\ 00000\ 00000$
$n = 3$	
$0.00000\ 00000\ 00000$	$0.88888\ 88888\ 88889$
$\pm 0.77459\ 66692\ 41483$	$0.55555\ 55555\ 55556$
$n = 4$	
$\pm 0.33998\ 10435\ 84856$	$0.65214\ 51548\ 62546$
$\pm 0.86113\ 63115\ 94053$	$0.34785\ 48451\ 37454$
$n = 5$	
$0.00000\ 00000\ 00000$	$0.56888\ 88888\ 88889$
$\pm 0.53846\ 93101\ 05683$	$0.47862\ 86704\ 99366$
$\pm 0.90617\ 98459\ 38664$	$0.23692\ 68850\ 56189$

[35, 47] — с 20 значащими цифрами вплоть до $n = 96$ и $n = 48$. Таким образом, практическое применение квадратур Гаусса обычно не встречает затруднений.

При небольших значениях n можно дать точные аналитические выражения для узлов формул Гаусса, как корней полиномов Лежандра $L_n(x)$, имеющих явные представления (2.106). Так, для $n = 3$

$$L_3(x) = \frac{1}{2}(5x^3 - 3x) = \frac{1}{2}x(5x^2 - 3).$$

Поэтому для канонического интервала интегрирования $[-1, 1]$ и $n = 3$

узлы квадратурной формулы Гаусса суть

$$\begin{aligned}x_1 &= -\sqrt{\frac{3}{5}} = -0.77459\ 66692\ 41483\dots, \\x_2 &= 0, \\x_3 &= \sqrt{\frac{3}{5}} = 0.77459\ 66692\ 41483\dots\end{aligned}$$

Для $n = 4$

$$L_4(x) = \frac{1}{8}(35x^4 - 30x^2 + 3),$$

и нахождение корней этого биквадратного полинома труда не представляет. Аналогично и для $n = 5$, когда

$$L_5(x) = \frac{1}{8}(63x^5 - 70x^3 + 15x) = \frac{1}{8}x(63x^4 - 70x^2 + 15).$$

Соответствующие весовые коэффициенты можно легко найти решением небольших систем линейных уравнений, к которым редуцируется система (2.126) после подстановки в неё известных значений узлов.

Численные значения узлов и весов квадратурных формул Гаусса для $n = 2, 3, 4, 5$ сведены в Табл. 2.3. Видно, что узлы располагаются симметрично относительно середины интервала интегрирования, а равноотстоящие от неё весовые коэффициенты одинаковы. Симметрия расположения узлов очевидно следует из того, что любой полином Лежандра является, в зависимости от номера, либо чётной, либо нечётной функцией.

2.13д Погрешность квадратур Гаусса

Для исследования остаточного члена квадратурных формул Гаусса предположим, что подинтегральная функция $f(x)$ имеет достаточно высокую гладкость. Построим для неё интерполяционный многочлен, принимающий в узлах x_1, x_2, \dots, x_n значения $f(x_1), f(x_2), \dots, f(x_n)$. Коль скоро квадратурная формула Гаусса точна на полиномах степени $2n - 1$, то для адекватного учёта этого факта степень полинома, интерполирующего подинтегральную функцию, также нужно взять равной $2n - 1$. Необходимую степень можно получить, рассматривая, как и для формулы Симпсона, интерполяцию с кратными узлами, например, назначая кратность всех n узлов равной двум и предполагая в них известными фиктивные значения производных $f'(x_1), f'(x_2), \dots, f'(x_n)$.

При этом согласно (2.44) погрешность интерполирования функции $f(x)$ полиномом Эрмита $H_{2n-1}(x)$ равна

$$\begin{aligned} R_{2n-1}(f, x) &= f(x) - H_{2n-1}(x) \\ &= \frac{f^{(2n)}(\xi(x))}{(2n)!} \cdot \prod_{i=1}^n (x - x_i)^2, \\ &= \frac{f^{(2n)}(\xi(x))}{(2n)!} \cdot (\omega(x))^2, \end{aligned}$$

где $\omega(x) = (x - x_1)(x - x_2) \cdots (x - x_n)$ и $\xi(x)$ — некоторая точка, зависящая от x , из открытого интервала интерполирования $]a, b[$. По условиям интерполяции $H_{2n-1}(x_i) = f(x_i)$, $i = 1, 2, \dots, n$, следовательно,

$$\begin{aligned} \int_a^b f(x) dx &= \int_a^b (H_{2n-1}(x) + R_{2n-1}(f, x)) dx \\ &= \int_a^b H_{2n-1}(x) dx + \int_a^b R_{2n-1}(f, x) dx \\ &= \sum_{i=1}^n c_i H_{2n-1}(x_i) + \int_a^b R_{2n-1}(f, x) dx \\ &= \sum_{i=1}^n c_i f(x_i) + \frac{1}{(2n)!} \int_a^b f^{(2n)}(\xi(x)) (\omega(x))^2 dx, \end{aligned}$$

где c_i — веса квадратурной формулы Гаусса.

Выражение для второго слагаемого последней суммы, т. е. для остаточного члена квадратуры, можно упростить и далее, приняв во внимание знакопостоянство множителя $(\omega(x))^2$. Тогда в силу интегральной теоремы о среднем (см., к примеру, [34]) имеем

$$\int_a^b f^{(2n)}(\xi(x)) (\omega(x))^2 dx = f^{(2n)}(\theta) \int_a^b (\omega(x))^2 dx$$

для некоторой точки $\theta \in]a, b[$. Таким образом, погрешность квадратурной формулы Гаусса, построенной по n узлам $x_1, x_2, \dots, x_n \in [a, b]$,

равна

$$R(f) = \frac{f^{(2n)}(\theta)}{(2n)!} \int_a^b (\omega(x))^2 dx,$$

где $\theta \in]a, b[$.

Узлы x_1, x_2, \dots, x_n — это корни полинома, полученного из полинома Лежандра линейной заменой переменных. По этой причине интеграл в полученной формуле для погрешности можно вычислить точно, приведя его заменой переменных к интервалу $[-1, 1]$ и воспользовавшись результатом Предложения 2.11.1. Это даёт

$$R(f) = \frac{(n!)^4}{((2n)!)^3(2n+1)} (b-a)^{2n+1} f^{(2n)}(\theta) \quad (2.137)$$

для некоторой промежуточной точки θ из интервала интегрирования $]a, b[$. Иногда удобнее грубая оценка

$$|R(f)| \leq \frac{(n!)^4}{((2n)!)^3(2n+1)} M_{2n} (b-a)^{2n+1},$$

где, как обычно, обозначено $M_p = \max_{x \in [a, b]} |f^{(p)}(x)|$.

В частности, для формулы Гаусса (2.130)–(2.131) с двумя узлами

$$|R_2(f)| \leq \frac{M_4(b-a)^5}{4320},$$

что даже лучше оценки погрешности для формулы Симпсона. Практическое поведение этой погрешности мы могли видеть в Примере 2.13.1.

Отметим, что выведенная оценка (2.137) справедлива лишь при достаточной гладкости подинтегральной функции $f(x)$. Вообще, квадратурные формулы Гаусса с большим числом узлов целесообразно применять лишь для функций, обладающих значительной гладкостью.

Другое важное наблюдение состоит в том, что в выражении (2.137) числитель $(n!)^4$ с ростом n может быть сделан сколь угодно меньшим знаменателя, превосходящего $((2n)!)^4$. Как следствие, если производные подинтегральной функции не растут «слишком быстро» с увеличением порядка, то с ростом числа узлов и гладкости интегрируемой функции порядок точности квадратурных формул Гаусса может быть сделан сколь угодно высоким. В этом квадратуры Гаусса принципиально отличаются, к примеру, от интерполяции с помощью сплайнов, которая сталкивается с ограничением на порядок сходимости, не зависящим

от гладкости исходных данных (стр. 96). Таким образом, квадратурные формулы Гаусса дают пример *ненасыщаемого* численного метода, порядок точности которого может быть сделан любым в зависимости от того, насколько гладкими являются входные данные для этого метода.

В заключение темы следует сказать, что на практике нередко требуется включение во множество узлов квадратурной формулы каких-либо фиксированных точек интервала интегрирования, например, его концов (одного или обоих), либо каких-то выделенных внутренних точек. Основная идея формул Гаусса может быть применена и в этом случае, что приводит к *квадратурам Маркова* [3, 11, 53, 56], которые называются также *квадратурами Лобатто* [2, 35].

Построение квадратурных формул Гаусса основывалось на оптимизации алгебраической степени точности квадратур. Эта идея может быть модифицирована и приспособлена к другим ситуациям, когда точность результата для алгебраических полиномов уже не являются наиболее адекватным мерилем качества квадратурной формулы. Например, можно развивать квадратуры наивысшего *тригонометрического порядка точности*, которые окажутся практичнее при вычислении интегралов от осциллирующих функций [56].

2.14 Составные квадратурные формулы

Рассмотренные выше квадратурные формулы дают приемлемую погрешность в случае, когда длина интервала интегрирования $[a, b]$ невелика и подинтегральная функция имеет на нём ограниченные производные. Но если величина $(b - a)$ относительно велика или интегрируемая функция имеет большие производные, то погрешность вычисления интеграла делается значительной или даже большой. Тогда для получения требуемой точности вычисления интеграла применяют *составные квадратурные формулы*, основанные на разбиении интервала интегрирования на подинтервалы меньшей длины, по каждому из которых вычисляется значение «элементарной квадратуры», а затем искомый интеграл приближается их суммой.¹⁸

Зафиксируем некоторую квадратурную формулу. Будем также считать, что её погрешность $R(f)$ имеет оценку

$$|R(f)| \leq C(b - a)^p,$$

¹⁸Интересно, что при этом подинтегральная функция интерполируется на всём интервале интегрирования $[a, b]$ при помощи сплайна.

где C — константа, зависящая от типа квадратурной формулы и интегрируемой функции, но в любом случае $p > 1$. К примеру, для формулы средних прямоугольников $C = M_2/24$ и $p = 3$, а для формулы Симпсона $C = M_4/2880$ и $p = 5$. Разобьём интервал интегрирования точками r_1, r_2, \dots, r_{N-1} на N равных частей $[a, r_1], [r_1, r_2], \dots, [r_{N-1}, b]$ длины $h = (b - a)/N$. Тогда, пользуясь аддитивностью интеграла, можно вычислить по рассматриваемой формуле интегралы

$$\int_{r_i}^{r_{i+1}} f(x) dx, \quad i = 0, 1, \dots, N - 1,$$

где $r_0 = a$, $b = r_N$, а затем положить

$$\int_a^b f(x) dx \approx \sum_{i=0}^{N-1} \int_{r_i}^{r_{i+1}} f(x) dx. \quad (2.138)$$

На каждом подинтервале $[r_i, r_{i+1}]$

$$|R(f)| \leq C \left(\frac{b-a}{N} \right)^p = Ch^p,$$

а полная погрешность интегрирования $\tilde{R}(f)$ при использовании представления (2.138) не превосходит суммы погрешностей отдельных слагаемых, т. е.

$$|\tilde{R}(f)| \leq N C \left(\frac{b-a}{N} \right)^p = \frac{C(b-a)^p}{N^{p-1}} = C(b-a)h^{p-1}.$$

Как видим, эта погрешность уменьшилась в N^{p-1} раз, и потенциально таким способом погрешность вычисления интеграла можно сделать сколь угодно малой.

Число $(p-1)$ часто называют *порядком точности* (составной) квадратурной формулы, и это понятие очевидно согласуется с данным ранее Определением 2.8.1. Ясно, что основная идея составных квадратурных формул работает и в случае неравномерного разбиения интервала интегрирования на более мелкие части, но анализ погрешности проводить тогда труднее.

Для равномерного разбиения интервала интегрирования составные квадратурные формулы выглядят особенно просто. Выпишем их явный вид для рассмотренных выше простейших квадратур Ньютона-Котеса и разбиения интервала интегрирования $[a, b]$ на N равных частей $[r_0, r_1], [r_1, r_2], \dots, [r_{N-1}, r_N]$ длины $h = (b - a)/N$ каждая, в котором $a = r_0$ и $r_N = b$.

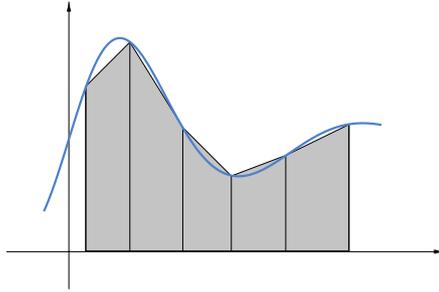


Рис. 2.24. Составная квадратурная формула трапеций

Составная формула (средних) прямоугольников

$$\int_a^b f(x) dx \approx h \sum_{i=1}^N f(r_{i-1/2}),$$

где $r_{i-1/2} = r_i - h/2$. Её полная погрешность

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

т.е. она имеет второй порядок точности. Эта формула, как нетрудно видеть, совпадает с интегральной суммой Римана для интеграла от $f(x)$ по интервалу $[a, b]$.

Составная формула трапеций

$$\int_a^b f(x) dx \approx h \left(\frac{1}{2}f(a) + \sum_{i=1}^{N-1} f(r_i) + \frac{1}{2}f(b) \right).$$

Её полная погрешность

$$|\tilde{R}(f)| \leq M_2 \frac{(b-a)h^2}{12},$$

т.е. порядок точности тоже второй.

Составная формула Симпсона (парабол)

$$\int_a^b f(x) dx \approx \frac{h}{6} \sum_{i=1}^N (f(r_{i-1}) + 4f(r_{i-1/2}) + f(r_i)),$$

где $r_{i-1/2} = r_i - h/2$. Её полная погрешность

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

т. е. формула имеет четвёртый порядок точности.

Аналогично можно сконструировать составные квадратурные формулы Гаусса, но мы не будем здесь развёртывать детали этого построения.

В составных квадратурных формулах увеличение точности вычисления интеграла достигается ценой дополнительных трудозатрат. В рассмотренном нами одномерном случае эти трудозатраты растут всего лишь линейно, хотя и здесь необходимость вычисления сложной подинтегральной функции может иногда быть весьма обременительной. Но при возрастании размерности интеграла, когда необходимо прибегнуть к составным кубатурным формулам, рост трудозатрат делается уже значительным, имея тот же порядок, что и размерность пространства. Так же растёт и погрешность суммирования результатов интегрирования по отдельным подобластям общей области интегрирования. Поэтому эффект увеличения точности составной формулы при возрастании размерности становится всё менее ощутимым.

2.15 Сходимость квадратур

С теоретической точки зрения интересен вопрос о сходимости квадратур при неограниченном возрастании числа узлов. Иными словами, верно ли, что

$$\sum_{k=0}^n c_k f(x_k) \rightarrow \int_a^b f(x) dx$$

при $n \rightarrow \infty$ (здесь узлы и веса квадратурных формул нумеруются с нуля)?

Похожий вопрос вставал при исследовании интерполяционного процесса, и мы обсуждали его в §2.5. Но в случае квадратурных формул

помимо бесконечной треугольной матрицы узлов

$$\begin{pmatrix} x_0^{(1)} & & & \cdots \\ x_0^{(2)} & x_1^{(2)} & & \mathbf{0} & \cdots \\ x_0^{(3)} & x_1^{(3)} & x_2^{(3)} & & \cdots \\ \vdots & \vdots & \vdots & \ddots & \ddots \end{pmatrix}, \quad (2.139)$$

таких что $x_k^{(n)}$ лежат на интервале интегрирования $[a, b]$ и $x_i^{(n)} \neq x_j^{(n)}$ при $i \neq j$, необходимо задавать ещё и треугольную матрицу весовых коэффициентов квадратурных формул

$$\begin{pmatrix} c_0^{(1)} & & & \cdots \\ c_0^{(2)} & c_1^{(2)} & & \mathbf{0} & \cdots \\ c_0^{(3)} & c_1^{(3)} & c_2^{(3)} & & \cdots \\ \vdots & \vdots & \vdots & \ddots & \ddots \end{pmatrix}. \quad (2.140)$$

В случае задания бесконечных треугольных матриц (2.139)–(2.140), по которым организуется приближённое вычисление интегралов на последовательности сеток, будем говорить, что определён *квадратурный процесс*.

Определение 2.15.1 *Квадратурный процесс, задаваемый зависящим от целочисленного параметра n семейством квадратурных формул*

$$\int_a^b f(x) dx \approx \sum_{k=0}^n c_k^{(n)} f(x_k^{(n)}), \quad n = 0, 1, 2, \dots,$$

которые определяются матрицами узлов и весов (2.139)–(2.140), будем называть *сходящимся для функции $f(x)$ на интервале $[a, b]$, если*

$$\lim_{n \rightarrow \infty} \sum_{k=0}^n c_k^{(n)} f(x_k^{(n)}) = \int_a^b f(x) dx,$$

т. е. если при неограниченном возрастании числа узлов n предел результатов квадратурных формул равен точному интегралу от функции f по $[a, b]$.

Необходимо оговориться, что в практическом плане вопрос о сходимости квадратур решается положительно с помощью составных формул, рассмотренных в предыдущем §2.14. Для достаточно общих функций путём построения составной квадратурной формулы всегда можно добиться сходимости приближённого значения интеграла к точному (для составной формулы прямоугольников это следует из самого определения интегрируемости по Риману). Обсуждаемый ниже круг вопросов относится больше к теоретическим качествам тех или иных «чистых» квадратурных формул, их предельному поведению.

Весьма общие достаточные условия для сходимости квадратур были сформулированы и обоснованы В.А.Стекловым [62], а впоследствии Д. По́йа [72] доказал также необходимость условий Стеклова.

Теорема 2.15.1 (теорема Стеклова-По́йа) *Квадратурный процесс, порождаемый матрицами узлов и весов (2.139)–(2.140), сходится для любой непрерывной на $[a, b]$ функции тогда и только тогда, когда*

- (1) *этот процесс сходится для полиномов,*
- (2) *суммы абсолютных значений весов равномерно по n ограничены, т. е. существует такая константа C , что*

$$\sum_{k=0}^n |c_k^{(n)}| \leq C \quad (2.141)$$

для всех $n = 0, 1, 2, \dots$

Покажем достаточность условий теоремы Стеклова-По́йа. С этой целью, задавшись каким-то $\epsilon > 0$, найдём полином $P_N(x)$, который равномерно с погрешностью ϵ приближает непрерывную подинтегральную функцию $f(x)$ на рассматриваемом интервале $[a, b]$. Существование такого полинома обеспечивается теоремой Вейерштрасса (см. §2.5). Далее преобразуем выражение для остаточного члена квадратурной форму-

лы:

$$\begin{aligned}
 R_n(f) &= \int_a^b f(x) \, dx - \sum_{k=0}^n c_k^{(n)} f(x_k^{(n)}) \\
 &= \int_a^b (f(x) - P_N(x)) \, dx + \int_a^b P_N(x) \, dx - \sum_{k=0}^n c_k^{(n)} f(x_k^{(n)}) \\
 &= \int_a^b (f(x) - P_N(x)) \, dx \\
 &\quad + \left(\int_a^b P_N(x) \, dx - \sum_{k=0}^n c_k^{(n)} P_N(x_k^{(n)}) \right) \\
 &\quad + \sum_{k=0}^n c_k^{(n)} (P_N(x_k^{(n)}) - f(x_k^{(n)})).
 \end{aligned}$$

Отдельные слагаемые полученной суммы, расположенные выше в различных строчках, оцениваются при достаточно больших номерах n следующим образом:

$$\left| \int_a^b (f(x) - P_N(x)) \, dx \right| \leq \epsilon (b - a), \quad \text{так как } P_N(x) \text{ приближает } f(x) \text{ равномерно с погрешностью } \epsilon \text{ на интервале } [a, b],$$

$$\left| \int_a^b P_N(x) \, dx - \sum_{k=0}^n c_k^{(n)} P_N(x_k^{(n)}) \right| \leq \epsilon, \quad \text{так как квадратуры сходятся на полиномах,}$$

$$\left| \sum_{k=0}^n c_k^{(n)} (P_N(x_k^{(n)}) - f(x_k^{(n)})) \right| \leq \epsilon \sum_{k=0}^n |c_k^{(n)}| \leq \epsilon C \quad \text{в силу (2.141).}$$

Поэтому в целом, если n достаточно велико, имеем

$$|R_n(x)| \leq \epsilon (b - a + 1 + C).$$

Это и означает сходимость рассматриваемого квадратурного процесса.

Доказательство необходимости условия теоремы Стеклова-Пойа помимо оригинальной статьи [72] можно найти также в книгах [3, 25].

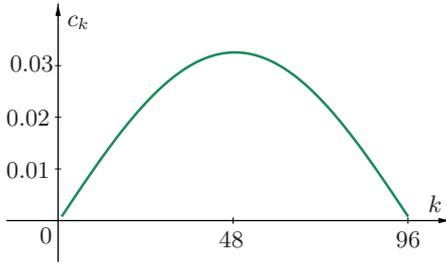


Рис. 2.25. Зависимость весовых коэффициентов от номера для квадратуры Гаусса 96-го порядка

В формулировке теоремы фигурирует величина

$$\sum_{k=0}^n |c_k|, \quad (2.142)$$

— сумма абсолютных значений весов, которая, как мы видели в §2.12а, является коэффициентом увеличения погрешности в данных и играет очень важную роль при оценке качества различных квадратурных формул. В §2.12д уже упоминался результат Р.О. Кузьмина [48] о том, что для формул Ньютона-Котеса величина (2.142) неограниченно увеличивается с ростом числа узлов n . Как следствие, на произвольных непрерывных функциях эти квадратурные формулы сходимостью не обладают.

Для квадратурных формул Гаусса ситуация другая. Справедливо

Предложение 2.15.1 *Весовые коэффициенты квадратурных формул Гаусса положительны.*

Доказательство. Ранее мы уже выводили для весов интерполяционных квадратурных формул выражение (2.122). Зафиксировав индекс $i \in \{1, 2, \dots, n\}$, дадим другое явное представление для весового коэффициента c_i квадратурной формулы Гаусса, из которого и будет следовать доказываемое предложение.

Пусть x_1, x_2, \dots, x_n — узлы квадратурной формулы Гаусса на интервале интегрирования $[a, b]$. Так как формулы Гаусса имеют алгебраическую степень точности $2n - 1$, то для полинома

$$P_i(x) = ((x - x_1) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n))^2$$

степени $2(n - 1)$ должно выполняться точное равенство

$$\int_a^b P_i(x) dx = \sum_{k=1}^n c_k P_i(x_k). \quad (2.143)$$

Но $P_i(x_k) = \delta_{ik}$ по построению полинома P_i , так что от суммы справа в (2.143) остаётся лишь одно слагаемое $c_i P_i(x_i)$:

$$\int_a^b P_i(x) dx = c_i P_i(x_i).$$

Следовательно,

$$c_i = \int_a^b P_i(x) dx / P_i(x_i).$$

Далее, $P_i(x) > 0$ всюду на интервале интегрирования $[a, b]$ за исключением конечного числа точек, и потому положителен интеграл в числителе выписанного выражения. Кроме того, $P_i(x_i) > 0$, откуда можно заключить, что $c_i > 0$. ■

Напомним, что сумма весов формул Гаусса равна длине интервала интегрирования (как и для всех интерполяционных квадратурных формул, см. §2.12г). Как следствие, величина (2.142) при этом ограничена, и квадратурный процесс по формулам Гаусса всегда сходится.

Завершая тему, можно отметить, что ситуация со сходимостью квадратур оказывается в целом более благоприятной, чем для интерполяционных процессов.

2.16 Вычисление интегралов методом Монте-Карло

В *методе Монте-Карло*, называемом также *методом статистического моделирования*, искомое решение задачи представляется в виде какой-либо характеристики специально построенного случайного процесса. Затем этот процесс моделируется, с помощью ЭВМ или какими-то другими средствами, и по его реализации мы вычисляем нужную характеристику, т. е. решение задачи. Наиболее часто решение задач представляется так называемым математическим ожиданием (средним значением) специально подобранной случайной величины.

В качестве примера рассмотрим задачу вычисления определённого интеграла

$$\int_a^b f(x) dx \quad (2.144)$$

от непрерывной функции $f(x)$. Согласно известной из интегрального исчисления теореме о среднем (см., к примеру, [34])

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

для некоторой точки $c \in [a, b]$. Смысл «средней точки» c можно понять глубже с помощью следующего рассуждения. Пусть интервал интегрирования $[a, b]$ разбит на N равных подинтервалов. По определению интеграла Римана, если x_i — точки из этих подинтервалов, то

$$\int_a^b f(x) dx \approx \sum_{i=1}^N \frac{b-a}{N} f(x_i) = (b-a) \cdot \frac{1}{N} \sum_{i=1}^N f(x_i)$$

для достаточно больших N . Сумма в правой части — это произведение ширины интервала интегрирования $(b-a)$ на среднее арифметическое значений подинтегральной функции f в точках x_i , $i = 1, 2, \dots, N$. Таким образом, интеграл от $f(x)$ по $[a, b]$ есть не что иное, как «среднее значение» функции $f(x)$ на интервале $[a, b]$, умноженное на ширину этого интервала.

Но при таком взгляде на искомый интеграл нетрудно заметить, что «среднее значение» функции $f(x)$ можно получить каким-либо существенно более эффективным способом, чем простое увеличение количества равномерно расположенных точек x_i . Например, можно попытаться раскидывать эти точки случайно по $[a, b]$, но «приблизительно равномерно». Резон в таком образе действий следующий: случайный, но равномерный выбор точек x_i позволит в пределе иметь то же «среднее значение» функции, но, возможно, полученное быстрее, так как при случайном бросании есть надежда, что будут легче учтены почти все «представительные» значения функции на $[a, b]$.

Для формализации высказанных идей целесообразно привлечь аппарат теории вероятностей. Эта математическая дисциплина исследует случайные явления, которые подчиняются свойству «статистической устойчивости», т. е. обнаруживают закономерности поведения в больших сериях повторяющихся испытаний. Одними из основных понятий

теории вероятностей являются понятия *вероятности*, *случайной величины* и её *функции распределения*. Случайной величиной называется переменная величина, значения которой зависят от случая и для которой определена так называемая функция распределения вероятностей. Вероятность — это величина, выражающая относительную частоту интересующего нас события, которая обычно устанавливается в большой серии испытаний. Функция распределения показывает, следовательно, вероятность появления тех или иных значений этой случайной величины. Конкретное значение, которое случайная величина принимает в результате отдельного опыта, обычно называют реализацией случайной величины.

Случайные и «приблизительно равномерные» точки моделируются так называемым равномерным вероятностным распределением, в котором при большом количестве испытаний (реализаций) в любые подинтервалы исходного интервала $[a, b]$, имеющие равную длину, попадает примерно одинаковое количество точек.

На этом пути мы и приходим к простейшему методу Монте-Карло для вычисления определённого интеграла (2.144):

<p>фиксируем натуральное число N ;</p> <p>организуем реализации $\xi_i, i = 1, 2, \dots, N$, для случайной величины ξ, имеющей равномерное распределение на интервале $[a, b]$;</p> <p>(искомый интеграл) $\leftarrow \frac{b-a}{N} \cdot \sum_{i=1}^N f(\xi_i)$;</p>	(2.145)
--	---------

Получение равномерно распределённой случайной величины (как и других случайных распределений) является не вполне тривиальной задачей. Но она удовлетворительно решена на существующем уровне развития вычислительной техники и информатики. Так, практически во всех современных языках программирования имеются средства для моделирования простейших случайных величин, в частности, равномерного распределения на интервале.

Рассмотрим теперь задачу определения площади фигуры с криволинейными границами (Рис. 2.26). Погрузим её в прямоугольник со

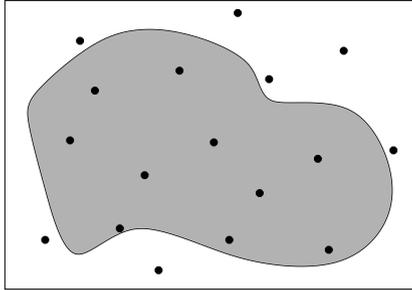


Рис. 2.26. Вычисление объёма области методом Монте-Карло.

сторонами, параллельными координатным осям, имеющий известные размеры, и станем случайным образом раскидывать точки внутри этого прямоугольника. Ясно, что при равномерном распределении случайных бросаний вероятность попадания точки в рассматриваемую фигуру равна отношению площадей этой фигуры и объемлющего её прямоугольника. С другой стороны, это отношение будет приблизительно равно относительной доле количества точек, которые попали в фигуру. Оно может быть вычислено в достаточно длинной серии случайных бросаний точек в прямоугольник.

На основе сформулированной выше идеи можно реализовать ещё один способ вычисления интеграла от функции одной переменной. Помещаем криволинейную трапецию, ограниченную графиком интегрируемой функции, в прямоугольник на плоскости Oxy . Затем организуем равномерное случайное бросание точек в этом прямоугольнике и подсчитываем относительную частоту точек, попадающих ниже графика интегрируемой функции. Искомый интеграл равен её произведению на площадь большого прямоугольника (см. Рис. 2.27).

Результаты вычислений по методу Монте-Карло сами являются случайной величиной, и два результата различных решений одной и той же задачи, вообще говоря, могут отличаться друг от друга. Можно показать, что второй (геометрический) способ вычисления интеграла методом Монте-Карло, вообще говоря, уступает по качеству результатов первому способу, основанному на нахождении «среднего значения» функции, так как дисперсия (среднеквадратичный разброс) получаемых оценок у него больше [28].

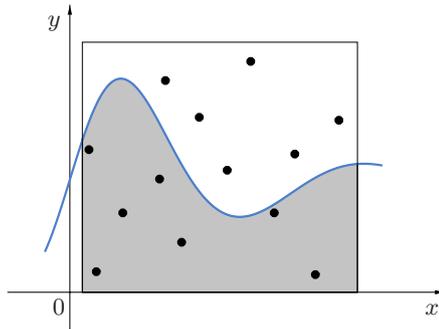


Рис. 2.27. Один из способов приближённого вычисления определённого интеграла методом Монте-Карло

Сформулированные выше идеи и основанные на них алгоритмы в действительности применимы для интегрирования функций от произвольного количества переменных. Более того, вероятностные оценки погрешности, пропорциональные $N^{-1/2}$, также не зависят от размерности n пространства, в котором берётся интеграл, тогда как для традиционных детерминистских методов интегрирования они ухудшаются с ростом n . Начиная с 7–8 переменных методы Монте-Карло уже превосходят по своей эффективности классические кубатурные формулы и являются главным методом вычисления многомерных интегралов.

В заключение параграфа — краткий исторический очерк. Идея моделирования случайных явлений очень стара. В современной истории науки использование статистического моделирования для решения конкретных практических задач можно отсчитывать с конца XVIII века, когда Ж.-Л. Бюффеном (в 1777 году) был предложен способ определения числа π с помощью случайных бросаний иглы на бумагу, разграфлённую параллельными линиями.¹⁹ Тем не менее, идея использования случайности при решении различных задач не получила большого развития вплоть до Второй мировой войны, т. е. до середины XX века.

В 1944 году в связи с работами по созданию атомной бомбы в США, поставившими ряд очень больших и сложных задач, С. Улам и Дж. фон Нейман предложили широко использовать для их решения статисти-

¹⁹Наиболее известная «докомпьютерная» реализация метода Бюффона была осуществлена американским астрономом А. Холлом [69].

ческое моделирование и аппарат теории вероятностей.²⁰ Этому способствовало появление к тому времени электронных вычислительных машин, позволивших быстро выполнять многократные статистические испытания (Дж. фон Нейман также принимал активное участие в создании первых цифровых ЭВМ). С конца 40-х годов XX века начинается широкое развитие метода Монте-Карло и методов статистического моделирования во всём мире. В настоящее время их успешно применяют для решения самых разнообразных задач практики (см., к примеру, [28, 55] и цитированную там литературу).

2.17 Правило Рунге для оценки погрешности

Предположим, что нам необходимо численно найти интеграл или производную функции, либо решение дифференциального или интегрального уравнения, т. е. решить какую-либо задачу, где фигурирует сетка на интервале вещественной оси или в пространстве бóльшего числа измерений. Пусть для решения этой задачи применяется численный метод порядка p , так что главный член его погрешности равен Ch^p , где h — шаг рассматриваемой сетки, а C — величина, напрямую от h не зависящая. Как правило, значение C не известно точно и его нахождение непосредственно из исходных данных задачи является делом трудным и малоперспективным. Мы могли видеть, к примеру, что для задач интерполирования и численного интегрирования выражение для этой константы вовлекает оценки для производных высоких порядков от рассматриваемой функции либо её разделённые разности. Во многих случаях их практическое вычисление не представляется возможным, так что оценки эти носят, главным образом, теоретический характер. Аналогична ситуация и с другими постановками задач и погрешностями их решения.

К. Рунге принадлежит идея использовать для определения значения константы C результаты нескольких расчётов на различных сетках. Далее, после того как величина C будет определена, мы можем использовать её значение для практического оценивания погрешности приближённых решений нашей задачи, которые получаются с помощью выбранного численного метода.

²⁰Интересно, что примерно в те же самые годы в СССР решение аналогичных задач советского атомного проекта было успешно выполнено другими методами.

Предположим для простоты анализа, что численные решения рассматриваемой задачи рассчитаны на сетках с шагом h и $h/2$ и равны соответственно J_h и $J_{h/2}$, а точное решение есть J . Тогда

$$J_h - J \approx Ch^p,$$

$$J_{h/2} - J \approx C \left(\frac{h}{2}\right)^p = C \frac{h^p}{2^p}.$$

Вычитая второе равенство из первого, получим

$$J_h - J_{h/2} \approx Ch^p - C \frac{h^p}{2^p} = Ch^p \frac{2^p - 1}{2^p},$$

так что

$$C \approx \frac{2^p}{2^p - 1} \cdot \frac{J_h - J_{h/2}}{h^p}.$$

Зная константу C , можно уже находить оценку погрешности рассчитанных решений J_h , $J_{h/2}$ или любых других

Правило Рунге работает плохо, если главный член погрешности Ch^p не доминирует над последующими членами её разложения, которые соответствуют $(p+1)$ -ой и более высоким степеням шага сетки h . Это происходит, как правило, для сильно меняющихся решений.

Литература к главе 2

Основная

- [1] БАРАХНИН В.Б., ШАПЕЕВ В.П. *Введение в численный анализ*. – Санкт-Петербург–Москва–Краснодар: Лань, 2005.
- [2] БАХВАЛОВ Н.С., ЖИДКОВ Н.П., КОВЕЛЬКОВ Г.М. *Численные методы*. – Москва: Бином, 2003, а также другие издания этой книги.
- [3] БЕРЕЗИН И.С., ЖИДКОВ Н.П. *Методы вычислений. Т. 1–2*. – Москва: Наука, 1966.
- [4] ВЕРЖБИЦКИЙ В.М. *Численные методы. Части 1–2*. – Москва: «Оникс 21 век», 2005.
- [5] ВОЛКОВ Е.А. *Численные методы*. – Москва: Наука, 1987.
- [6] ГЕЛЬФОНД А.О. *Исчисление конечных разностей*. – Москва: Наука, 1967.
- [7] ГОНЧАРОВ В.Л. *Теория интерполирования и приближения функций*. – Москва: ГИТТЛ, 1954.
- [8] ДАУГАВЕТ И.К. *Введение в теорию приближения функций*. – Ленинград: Издательство Ленинградского университета, 1977.

- [9] Демидович Б.П., Марон А.А. *Основы вычислительной математики*. – Москва: Наука, 1970.
- [10] Завьялов Ю.С., Квасов Б.И., Мирошниченко В.Л. *Методы сплайн-функций*. – Москва: Наука, 1980.
- [11] Калиткин Н.Н. *Численные методы*. – Москва: Наука, 1978.
- [12] Кобков В.В., Шокин Ю.И. *Сплайн-функции в численном анализе*. – Новосибирск: Издательство НГУ, 1983.
- [13] Коллатц Л. *Функциональный анализ и вычислительная математика*. – Москва: Мир, 1969.
- [14] Кострикин А.Н. *Введение в алгебру. Часть 1. Основы алгебры*. – Москва: Физматлит, 2001.
- [15] Крылов А.Н. *Лекции о приближённых вычислениях*. – Москва: ГИТТЛ, 1954, а также более ранние издания.
- [16] Крылов В.И. *Приближённое вычисление интегралов*. – Москва: Наука, 1967.
- [17] Крылов В.И., Бобков В.В., Монастырный П.И. *Вычислительные методы. Т. 1–2*. – Москва: Наука, 1976.
- [18] Кунц К.С. *Численный анализ*. – Киев: Техника, 1964.
- [19] Люстерник Л.А., Червоненкис О.А., Янпольский А.Р. *Математический анализ. Вычисление элементарных функций*. – Москва: ГИФМЛ, 1963.
- [20] Мак-Кракен Д., Дорн У. *Численные методы и программирование на ФОР-ТРАНе*. – Москва: Мир, 1977.
- [21] Марков А.А. *Исчисление конечных разностей*. – Одесса: Mathesis, 1910.
- [22] Мацокин А.М., Сорокин С.Б. *Численные методы. Часть 1. Численный анализ*. – Новосибирск: НГУ, 2006.
- [23] Миньков С.Л., Миньков Л.Л. *Основы численных методов*. – Томск: Издательство научно-технической литературы, 2005.
- [24] Мысовских И.П. *Интерполяционные кубатурные формулы*. – Москва: Наука, 1981.
- [25] Натансон И.П. *Конструктивная теория функций*. – Москва–Ленинград: ГИТТЛ, 1949.
- [26] Никольский С.М. *Квадратурные формулы*. – Москва: Наука, 1988.
- [27] Самарский А.А., Гулин А.В. *Численные методы*. – Москва: Наука, 1989.
- [28] Соболев И.М. *Численные методы Монте-Карло*. – Москва: Наука, 1973.
- [29] Стечкин С.Б., Субботин Ю.Н. *Сплайны в вычислительной математике*. – Москва: Наука, 1976.
- [30] Тихонов А.Н., Арсенин В.Я. *Методы решения некорректных задач*. – Москва: Наука, 1974.
- [31] Тыртышников Е.Е. *Матричный анализ и линейная алгебра*. – Москва: Физматлит, 2007.

- [32] Тыртышников Е.Е. *Методы численного анализа*. – Москва: Академия, 2007.
- [33] УИТТЕКЕР Э., РОБИНСОН Г. *Математическая обработка результатов наблюдений*. – Ленинград-Москва: ГТТИ, 1933.
- [34] ФИХТЕНГОЛЬЦ Г.М. *Курс дифференциального и интегрального исчисления*. Т. 1, 2. – Москва: Наука, 1966.

Дополнительная

- [35] АБРАМОВИЦ М., СТИГАН И. *Таблицы специальных функций*. – Москва: Наука, 1979.
- [36] АЛБЕРГ Дж., НИЛЬСОН Э., УОЛШ Дж. *Теория сплайнов и её приложения*. – Москва: Мир, 1972.
- [37] БАБЕНКО К.И. *Основы численного анализа*. – Москва: Наука, 1986.
- [38] БАХВАЛОВ Н.С., КОРНЕВ А.А., ЧИЖОНКОВ Е.В. *Численные методы. Решение задач и упражнения*. – Москва: Дрофа, 2008.
- [39] ГЕРОНИМУС Я.Л. *Теория ортогональных многочленов*. – Москва: Госуд. изд-во технико-теоретической литературы, 1950.
- [40] ГУРВИЦ А., КУРАНТ Р. *Теория функций*. – Москва: Наука, Физматлит, 1968.
- [41] ДЕМИДЕНКО Е.З. *Оптимизация и регрессия*. – Москва: Наука, 1989.
- [42] ДРОБЫШЕВИЧ В.И., ДЫМНИКОВ В.П., РИВИН Г.С. *Задачи по вычислительной математике*. – Москва: Наука, 1980.
- [43] КАХАНЕР Д., МОУЛЕР К., НЭШ С. *Численные методы и программное обеспечение*. – Москва: Мир, 1998.
- [44] КВАСОВ Б.И. *Методы изогометрической аппроксимации сплайнами*. – Москва: Физматлит, 2006.
- [45] КОЛЛАТЦ Л., КРАВС В. *Теория приближений. Чебышёвские приближения и их приложения*. – Москва: Наука, 1978.
- [46] КРОНРОД А.С. *Узлы и веса квадратурных формул. Шестнадцатизначные таблицы*. – Москва: Наука, 1964.
- [47] КРЫЛОВ В.И., ШУЛЬГИНА Л.Т. *Справочная книга по численному интегрированию*. – Москва: Наука, 1966.
- [48] КУЗЬМИН Р.О. К теории механических квадратур // *Известия Ленинградского политехнического института. Отделение техн. естеств. и матем.* 1931. – Т. 33. – С. 5–14.
- [49] ЛИННИК Ю.В. *Метод наименьших квадратов и основы теории обработки наблюдений*. 2-е изд. – Москва: ГИФМЛ, 1962.
- [50] ЛОКУЦИЕВСКИЙ О.В., ГАВРИКОВ М.Б. *Начала численного анализа*. – Москва: ТОО «Янус», 1994.
- [51] ЛОРАН Ж.-П. *Аппроксимация и оптимизация*. – Москва: Мир, 1975.
- [52] МЕНЬШИКОВ Г.Г. *Локализуемые вычисления. Конспект лекций*. – Санкт-Петербург: СПбГУ, Факультет прикладной математики–процессов управления, 2003.

- [53] МИКЕЛАДЗЕ Ш.Е. *Численные методы математического анализа*. – Москва: ГИТТЛ, 1953.
- [54] МИЛН В.Э. *Численный анализ*. – Москва: Издательство иностранной литературы, 1951.
- [55] МИХАЙЛОВ Г.А., ВОЙТИШЕК А.В. *Численное статистическое моделирование. Методы Монте-Карло*. – Москва: Изд. центр «Академия», 2006.
- [56] МЫСОВСКИХ И.П. *Лекции по методам вычислений*. – Санкт-Петербург: Издательство Санкт-Петербургского университета, 1998.
- [57] ПАШКОВСКИЙ С. *Вычислительные применения многочленов и рядов Чебышёва*. – Москва: Наука, 1983.
- [58] ПОГОРЕЛОВ А.И. *Дифференциальная геометрия*. – Москва: Наука, 1974.
- [59] РЕМЕЗ Е.Я. *Основы численных методов чебышёвского приближения*. – Киев: Наукова думка, 1969.
- [60] СЕГЁ Г. *Ортогональные многочлены*. – Москва: Физматлит, 1962.
- [61] СОВОЛЕВ С.Л. *Введение в теорию кубатурных формул*. – Москва: Наука, 1974.
- [62] СТЕКЛОВ В.А. О приближённом вычислении определённых интегралов // *Известия Академии Наук*. – 1916. – Т. 10, №6. – С. 169–186.
- [63] СУЕТИН П.К. *Классические ортогональные многочлены*. – Москва: Наука, 1979.
- [64] СТЕФЕНСЕН И.Ф. *Теория интерполяции*. – Москва: Объединённое научно-техническое издательство НКТП СССР, 1935.
- [65] ХАНСЕН Э., УОЛСТЕР ДЖ.У. *Глобальная оптимизация с помощью методов интервального анализа*. – Москва-Ижевск: Издательство «РХД», 2012.
- [66] ХАУСХОЛДЕР А.С. *Основы численного анализа*. – Москва: Издательство иностранной литературы, 1956.
- [67] ХЕММИНГ Р.В. *Численные методы*. – Москва: Наука, 1972.
- [68] АВЕРТН О. *Precise numerical methods using C++*. – San Diego: Academic Press, 1998.
- [69] HALL A. On an experimental determination of π // *Messenger of Mathematics*. – 1873. – Vol. 2. – P. 113–114.
- [70] ЛОВАЧЕВСКИЙ Н. Probabilité des résultats moyens tirés d'observations répétées // *Journal für die reine und angewandte Mathematik*. – 1842. – Bd. 24. – S. 164–170.
- [71] MOORE R.E., KEARFOTT R.B., CLOUD M. *Introduction to interval analysis*. – Philadelphia: SIAM, 2009.
- [72] POLYA G. Über Konvergenz von Quadraturverfahren // *Mathematische Zeitschrift*. – 1933. – Bd. 37. – S. 264–286.
- [73] SCHOENBERG I.J Contributions to the problem of approximation of equidistant data by analytic functions. Part A: On the problem of smoothing or graduation. A first class of analytic approximation formulae. Part B: On the problem of osculatory interpolation. A second class of analytic approximation formulae // *Quart. Appl. Math.* – 1946. – Vol. 4. – P. 45–99, 112–141.

Глава 3

Численные методы линейной алгебры

3.1 Задачи вычислительной линейной алгебры

Численные методы линейной алгебры — это один из классических разделов вычислительной математики, который в середине XX века вычленился даже в отдельное научное направление в связи с бурным развитием математических вычислений на ЭВМ.¹ Традиционный, исторически сложившийся список задач вычислительной линейной алгебры по состоянию на 50–60-е годы прошлого века можно найти в капитальной книге Д.К. Фаддеева и В.Н. Фаддеевой [44]. Он включал

- решение систем линейных алгебраических уравнений,
- вычисление определителей матриц,
- нахождение обратной матрицы,
- нахождение собственных значений и собственных векторов матриц,

а также многочисленные разновидности этих задач.

¹В англоязычной учебной и научной литературе часто используют термин «матричные вычисления», который уже по объёму, не охватывая, к примеру, такую часть современной вычислительной линейной алгебры как тензорные вычисления.

Но «всё течёт, всё меняется». По мере развития науки и технологий в фокусе развития вычислительной линейной алгебры оказались новые задачи. Вот как формулирует список важнейших задач в 2001 году американский специалист Дж. Деммель в книге [13]:

- решение систем линейных алгебраических уравнений;
- линейная задача о наименьших квадратах:
найти вектор x , минимизирующий $\langle Ax - b, Ax - b \rangle$
для заданных $m \times n$ -матрицы A и m -вектора b ;
- нахождение собственных значений и собственных векторов матриц;
- нахождение сингулярных чисел и сингулярных векторов матриц.

Постановку последней задачи мы будем обсуждать ниже в §3.2г. Вторая задача из этого списка — линейная задача о наименьших квадратах — является одним из вариантов дискретной задачи о наилучшем среднеквадратичном приближении. Она возникает обычно в связи с решением переопределённых систем линейных алгебраических уравнений (СЛАУ), которые, к примеру, получаются при обработке экспериментальных данных.

Помимо перечисленных задач к сфере вычислительной линейной алгебры относится также решение разнообразных линейных уравнений, в которых неизвестными являются матрицы (матричные уравнения Сильвестера, Ляпунова и др.). Эти уравнения возникают, к примеру, в теории автоматического управления.

С точки зрения классических разделов математики решение выписанных задач даётся вполне конструктивными способами и как будто не встречает затруднений:

- решение квадратной СЛАУ получается покомпонентно по формуле Крамера, как частное двух определителей, которые, в свою очередь, могут быть вычислены по явной формуле;
- для вычисления собственных значений матрицы A нужно выписать её характеристическое (вековое) уравнение $\det(A - \lambda I) = 0$ и найти его корни λ .

И так далее. Но практическая реализация этих теоретических рецептов наталкивается на почти непреодолимые трудности.

К примеру, явная формула для определителя $n \times n$ -матрицы выражает его как сумму $n!$ слагаемых, каждое из которых есть произведение n элементов из разных строк и столбцов матрицы. Раскрытие определителя по этой формуле требует $n!(n-1)$ умножений и $(n!-1)$ сложений, т. е. всего примерно $n!n$ арифметических операций, и потому из-за взрывного роста факториала² решение СЛАУ по правилу Крамера при $n \approx 20-30$ делается невозможным даже на самых современных ЭВМ.

Производительность современных ЭВМ принято выражать в так называемых *флопах* (сокращение от английской фразы floating point operation), и 1 флоп — это одна усреднённая арифметическая операция в арифметике с плавающей точкой в секунду (см. §1.2). Для наиболее мощных на сегодняшний день ЭВМ скорость работы измеряется так называемым петафлопами, 10^{15} операций с плавающей точкой в секунду. Для круглого счёта можно даже взять производительность нашего гипотетического компьютера равной 1 эксафлоп = 10^{18} операций с плавающей точкой в секунду. Решение на такой вычислительной машине системы линейных алгебраических уравнений размера 30×30 по правилу Крамера, с раскрытием определителей по явной комбинаторной формуле, потребует времени

$$30 \text{ компонент решения} \cdot \frac{30 \cdot 30! \text{ операций}}{10^{18} \text{ флоп} \cdot 3600 \frac{\text{сек}}{\text{час}} \cdot 24 \frac{\text{час}}{\text{сутки}} \cdot 365 \frac{\text{сутки}}{\text{год}}}$$

лет, т. е. примерно $7.57 \cdot 10^9$ лет. Для сравнения, возраст Земли в настоящее время оценивается в $4.5 \cdot 10^9$ лет.

Обращаясь к задаче вычисления собственных значений матрицы, напомним известную из алгебры теорема Абеля-Руффини³: общее алгебраическое уравнение степени выше четвёртой «неразрешимо в радикалах», т. е. не существует конечная формула, выражающая решения такого уравнения через коэффициенты с помощью арифметических операций и взятия корней произвольной степени. Таким образом, для матриц размера 5×5 и более мы по необходимости должны развивать для нахождения собственных значений какие-то численные методы. К этому добавляются трудности в раскрытии определителя, который входит в характеристическое уравнение матрицы.

²Напомним в этой связи известную в математическом анализе асимптотическую формулу Стирлинга — $n! \approx \sqrt{2\pi n} (n/e)^n$, где $e = 2.7182818\dots$

³Иногда её называют просто «теоремой Абеля» (см., к примеру, [61]).

Помимо неприемлемой трудоёмкости ещё одной причиной непригодности для реальных вычислений некоторых широко известных алгоритмов из «чистой математики» является сильное влияние на их результаты неизбежных погрешностей счёта и ввода данных. Например, очень неустойчиво к ошибкам решение СЛАУ по правилу Крамера.

3.2 Теоретическое введение

3.2a Основные понятия

Термин «вектор» имеет несколько значений. Прежде всего, это направленный отрезок на прямой, плоскости или в пространстве. Далее, термин «вектор» может обозначать упорядоченный кортеж из чисел либо объектов какой-то другой природы, расположенный вертикально (вектор-столбец) или горизонтально (вектор-строка). Таким образом, если a_1, a_2, \dots, a_n — некоторые числа, то

$$a = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} \text{ — это вектор-столбец,}$$

а

$$a = (a_1, a_2, \dots, a_n) \text{ — это вектор-строка.}$$

Этот смысл термина «вектор» широко используется в информатике и программировании. Наконец, «вектор» можно понимать как элемент абстрактного «векторного пространства», и в современной математике огромное применение находят, к примеру, линейные векторные пространства, об элементах которых мы привычно говорим, как о некоторых «векторах».

Все три перечисленных выше смысла тесно связаны между собой и взаимно проникают друг в друга. Мы в равной степени будем пользоваться всеми ими, предполагая, что контекст изложения не даст повода к недоразумениям. По умолчанию, если не оговорено противное, условимся считать, что «векторами» во втором смысле являются вектор-столбцы. Множество векторов, компоненты которых принадлежат вещественной оси \mathbb{R} или комплексной плоскости \mathbb{C} , мы будем обозначать через \mathbb{R}^n или \mathbb{C}^n . При этом нулевые векторы, т. е. векторы, все компоненты которых суть нули, мы традиционно обозначаем через «0».

Ненулевые векторы a и b называются *коллинеарными*, если $a = \alpha b$ для некоторого скаляра α . Иногда различают *сонаправленные* коллинеарные векторы, отвечающие случаю $\alpha > 0$, и *противоположно направленные*, для которых $\alpha < 0$. Нулевой вектор по определению коллинеарен любому вектору.

Вообще, в линейной алгебре, при работе с линейными векторными пространствами, большую роль играют линейные выражения вида

$$\alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_k v_k,$$

где $\alpha_1, \alpha_2, \dots, \alpha_k$ — некоторые скаляры, а v_1, v_2, \dots, v_k — векторы из рассматриваемого пространства. Такие выражения называются *линейными комбинациями* векторов v_1, v_2, \dots, v_k . Говорят также, что линейная комбинация *нетривиальная*, если хотя бы один из коэффициентов $\alpha_1, \alpha_2, \dots, \alpha_k$ не равен нулю.

Линейной оболочкой векторов v_1, v_2, \dots, v_k называют множество всевозможных линейных комбинаций этих векторов, т.е. наименьшее линейное подпространство, содержащее эти векторы v_1, \dots, v_k . Мы будем обозначать линейную оболочку посредством $\text{lin}\{v_1, \dots, v_n\}$, так что

$$\text{lin}\{v_1, \dots, v_n\} = \left\{ \sum_{i=1}^n \alpha_i v_i \mid \alpha_i \in \mathbb{R} \right\}.$$

На линейных пространствах \mathbb{R}^n и \mathbb{C}^n можно задать *скалярные произведения*. Напомним, что в вещественном случае это положительно определённая симметричная и билинейная форма, а в комплексном — положительно определённая эрмитова форма. Обычно они задаются в следующем стандартном виде

$$\langle a, b \rangle = \sum_{i=1}^n a_i b_i \quad \text{для } a, b \in \mathbb{R}^n \quad (3.1)$$

или

$$\langle a, b \rangle = \sum_{i=1}^n a_i \bar{b}_i \quad \text{для } a, b \in \mathbb{C}^n. \quad (3.2)$$

Наличие скалярного произведения позволяет измерять углы между векторами и ввести очень важное понятие ортогональности векторов. При этом пространства \mathbb{R}^n и \mathbb{C}^n становятся евклидовыми пространствами, и для них справедливы многие красивые и важные свойства, существенно упрощающие математические рассуждения.

Матрица — это прямоугольная таблица, составленная из чисел или каких-нибудь других объектов. Если она имеет m строк и n столбцов, то обычно её записывают в виде

$$A := \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix},$$

называя a_{ij} *элементами* матрицы $A = (a_{ij})$. При этом мы можем отождествлять векторы с матрицами размера $n \times 1$ (вектор-столбцы) либо $1 \times n$ (вектор-строки).

Матрицы используются в современной математике для самых разнообразных целей. В частности, если дана система, составленная из конечного числа объектов (подсистем), то взаимодействие i -го объекта с j -ым можно описывать матрицей, элементы которой суть a_{ij} . В простейшем случае эти элементы принимают значения 1 или 0, соответствующие ситуациям «связь есть» и «никак не связано». Для нашего курса особенно важно, что с помощью матриц, как это показывается в линейной алгебре, даётся удобное представление для линейных отображений конечномерных пространств.

Ведущей подматрицей некоторой матрицы называется матрица, составленная из строк и столбцов с первыми номерами. *Ведущий минор* матрицы — это определитель ведущей подматрицы.

Транспонированной к $m \times n$ -матрице $A = (a_{ij})$ называется $n \times m$ -матрица A^T , в которой ij -ым элементом является a_{ji} . Иными словами,

$$A^T := \begin{pmatrix} a_{11} & a_{21} & \dots & a_{n1} \\ a_{12} & a_{22} & \dots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1m} & a_{2m} & \dots & a_{nm} \end{pmatrix}.$$

Для числовых матриц определены сумма, разность и произведение. Напомним, что сумма (разность) двух матриц одинакового размера есть матрица того же размера, образованная поэлементными суммами (разностями) операндов. Если $A = (a_{ij})$ — $m \times l$ -матрица и $B = (b_{ij})$ — $l \times n$ -матрица, то произведение матриц A и B есть такая $m \times n$ -матрица $C = (c_{ij})$, что

$$c_{ij} := \sum_{k=1}^l a_{ik} b_{kj}.$$

Строчным рангом матрицы (или рангом по строкам) называется, как известно, количество её линейно независимых строк. *Столбцовым рангом* матрицы (или рангом по столбцам) называется максимальное количество её линейно независимых столбцов. В курсах линейной алгебры показывается, что строчный и столбцовый ранги матрицы совпадают друг с другом и равны максимальному размеру ненулевого минора, порождённого этой матрицей. Как следствие, мы можем говорить просто о ранге матрицы.

Квадратная матрица, все строки которой (или столбцы) линейно независимы, называется *неособенной* (регулярной, невырожденной). Её ранг равен, таким образом, её порядку. В противном случае матрица называется *особенной* (вырожденной).

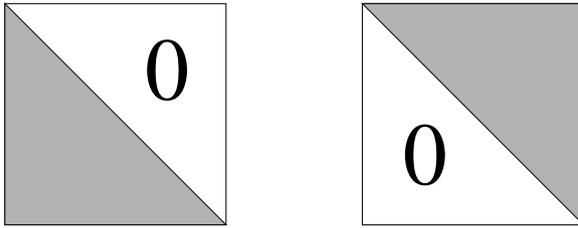


Рис. 3.1. Наглядные образы нижней треугольной и верхней треугольной матриц.

В случае, когда нулевые и ненулевые элементы в матрице A структурированы определённым образом, по отношению к A будут употребляться дополнительные определяющие термины. Например,

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ & a_{22} & \dots & a_{2n} \\ & & \ddots & \vdots \\ 0 & & & a_{nn} \end{pmatrix} \quad \text{и} \quad \begin{pmatrix} a_{11} & & & 0 \\ a_{21} & a_{22} & & \\ \vdots & \vdots & \ddots & \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}$$

— это *верхняя треугольная* и *нижняя треугольная* матрицы соответственно. Равносильные термины — *правая треугольная* и *левая треугольная* матрицы. Выбор того или иного варианта названия обычно диктуется контекстом или сложившейся традицией.

Обобщением понятия треугольных матриц на произвольный прямоугольный (неквадратный) случай являются *трапецеидальные матрицы*. Именно, прямоугольная матрица с нулями выше (ниже) диагонали называется нижней (верхней) трапецеидальной матрицей.

Блочными называются матрицы вида

$$\begin{pmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & A_{22} & \dots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \dots & A_{mn} \end{pmatrix},$$

элементы A_{ij} которых, в свою очередь, являются матрицами. Матрицы вида

$$\begin{pmatrix} A_{11} & & & 0 \\ & A_{22} & & \\ 0 & & \ddots & \\ & & & A_{nn} \end{pmatrix} \quad \text{и} \quad \begin{pmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ & A_{22} & \dots & A_{2n} \\ 0 & & \ddots & \vdots \\ & & & A_{nn} \end{pmatrix}$$

где внедиагональные блоки или же блоки ниже главной диагонали являются нулевыми, назовём соответственно *блочно-диагональными* или *верхними блочно треугольными* (правыми блочно треугольными), см. Рис. 3.2. Аналогичным образом определяются нижние блочно треугольные (левые блочно треугольные) матрицы.

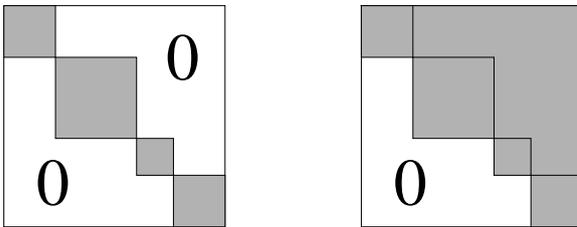


Рис. 3.2. Наглядные образы блочно-диагональной и верхней блочно-треугольной матриц.

Введение структурированных матриц и отдельное их изучение мотивируется тем, что многие операции с такими матрицами можно выполнить более специальным образом и существенно проще, чем в самом

общем случае. В частности, для блочных матриц операции выполняются «по блокам», т. е. совершенно аналогично операциям над обычными матрицами, но определённым «поблочным» образом, когда блоки выступают как отдельные самостоятельные элементы.

Линейная алгебра и её численные методы в некоторых ситуациях по существу требуют выхода в поле комплексных чисел \mathbb{C} , алгебраически пополняющее вещественную ось \mathbb{R} . Это необходимо, в частности, в связи с понятиями собственных чисел и собственных векторов матриц, но может также диктоваться исходной содержательной постановкой задачи. Например, привлечение комплексных чисел бывает необходимым при исследовании колебательных режимов в различных системах, так как в силу известной из математического анализа формулы Эйлера гармонические колебания с угловой частотой ω обычно представляются в виде комплексной экспоненты $\exp(i\omega t)$.

Эрмитово-сопряжённой к $m \times n$ -матрице $A = (a_{ij})$ называют $n \times m$ -матрицу A^* , в которой ij -ым элементом является комплексно-сопряжённый \bar{a}_{ji} . Иными словами,

$$A^* := \begin{pmatrix} \bar{a}_{11} & \bar{a}_{21} & \dots & \bar{a}_{n1} \\ \bar{a}_{12} & \bar{a}_{22} & \dots & \bar{a}_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ \bar{a}_{1m} & \bar{a}_{2m} & \dots & \bar{a}_{nm} \end{pmatrix},$$

и эрмитово сопряжение матрицы есть композиция транспонирования и взятия комплексного сопряжения элементов.

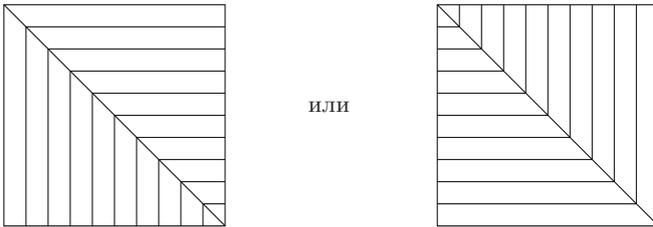


Рис. 3.3. Наглядные образы симметричной матрицы.

В линейной алгебре её приложениях широко используются специальные типы матриц — эрмитовы, симметричные, косоэрмитовы, косо-

симметричные, унитарные, ортогональные и т. п. Напомним, что *симметричными матрицами*⁴ называют матрицы, совпадающие со своими транспонированными, т. е. удовлетворяющие $A^T = A$. *Эрмитовыми матрицами* называются такие комплексные матрицы A , что $A^* = A$. Матрица Q называется *унитарной*, если $Q^*Q = I$. Матрица Q называется *ортогональной*, если $Q^T Q = I$.

Разреженными называются матрицы, большинство элементов которых равны нулю. Такие матрицы довольно часто встречаются в математическом моделировании, поскольку описывают системы или модели, в которых каждый элемент связан с относительно немногими другими элементами системы. Это происходит, например, если связи между элементами системы носят локальный характер. В противоположность этому, *плотно заполненными* называют матрицы, которые не являются разреженными. Иными словами, в плотно заполненных матрицах большинство элементов не равны нулю.

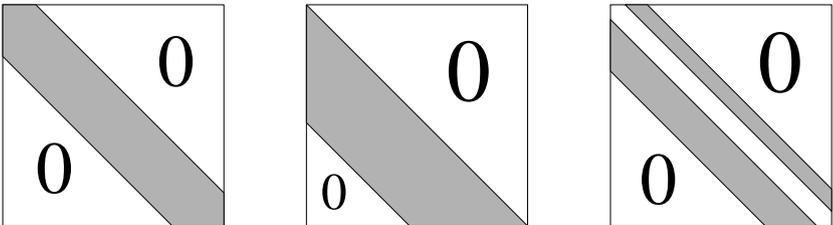


Рис. 3.4. Наглядные образы некоторых ленточных матриц.

Ленточными матрицами называют матрицы, у которых ненулевые элементы образуют выраженную «ленту» вокруг главной диагонали. В формальных терминах, матрица $A = (a_{ij})$ называется ленточной, если существуют такие натуральные числа p и q , что $a_{ij} = 0$ при $j - i > p$ и $i - j > q$. В этом случае величина $p + q + 1$ называется *шириной ленты*. Простейшими и важнейшими из ленточных матриц являются трёхдиагональные матрицы, для которых $p = q = 1$, и *двухдиагональные матрицы*, для которых $p = 0$ и $q = 1$ или $p = 1$ и $q = 0$. Такие матрицы встретятся нам в §3.8.

⁴Используют также термин *симметрическая матрица*.

3.2б Собственные числа и собственные векторы матрицы

Как должно быть известно читателю, для квадратных вещественных или комплексных матриц большую роль в теории и приложениях играют *собственные значения* и *собственные векторы*. Если обозначить посредством λ собственное значение $n \times n$ -матрицы A , а x , $x \neq 0$, — её собственный вектор, то они удовлетворяют матричному уравнению

$$Ax = \lambda x. \quad (3.3)$$

Содержательный смысл этого равенства состоит в том, что на одномерном линейном подпространстве в \mathbb{R}^n или \mathbb{C}^n , порождённом собственным вектором x , задаваемое матрицей A линейное преобразование действует как умножение на скаляр λ , т. е. как растяжение или сжатие. Собственные значения являются корнями так называемого *характеристического уравнения* матрицы, которое имеет вид $\det(A - \lambda I) = 0$. Для $n \times n$ -матрицы это алгебраическое уравнение n -ой степени, так что для его разрешимости по существу требуется привлечение поля комплексных чисел \mathbb{C} . Совокупность собственных чисел матрицы называется её *спектром*, и в общем случае спектр — подмножество комплексной плоскости.

Наконец, широко известный факт: собственные значения эрмитовых и симметричных матриц вещественны.

Предложение 3.2.1 Пусть A — $m \times n$ -матрица, B — $n \times m$ -матрица, так что одновременно определены произведения AB и BA . Спектры матриц AB и BA могут различаться только нулём.

Доказательство. Пусть λ — какое-нибудь ненулевое собственное значение матрицы AB , так что

$$ABu = \lambda u \quad (3.4)$$

с некоторым вектором $u \neq 0$. Умножая это равенство слева на матрицу B , получим

$$B(ABu) = B(\lambda u),$$

или

$$BA(Bu) = \lambda(Bu),$$

причём $Bv \neq 0$, так как иначе в исходном соотношении (3.4) необходимо должно быть $\lambda = 0$. Сказанное означает, что вектор Bv является собственным вектором матрицы BA , отвечающим такому же собственному значению λ .

И наоборот, если ненулевое μ есть собственное значение для BA , то, домножая слева равенство

$$BAv = \mu v$$

на матрицу A , получим

$$ABAv = AB(Av) = \mu(Av),$$

причём $Av \neq 0$. Как следствие, Av есть собственный вектор матрицы AB , отвечающий собственному значению μ . Иными словами, ненулевые собственные числа матриц AB и BA находятся во взаимнооднозначном соответствии друг с другом. ■

Другой вывод этого результата можно найти, к примеру, в [2, 42]. Особая роль нулевого собственного значения в этом результате объясняется тем, что если A и B — прямоугольные матрицы, то из двух матриц AB и BA по крайней мере одна имеет неполный ранг — та, чьи размеры больше. Но меньшая по размерам матрица может быть как особенной, так и неособенной.

Собственные векторы x , являющиеся решениями (3.3), называют также *правыми собственными векторами*, поскольку они соответствуют умножению на матрицу справа. Но нередко возникает необходимость рассмотрения *левых собственных векторов*, обладающих свойством, аналогичным (3.3), но при умножении на матрицу слева. Очевидно, это должны быть собственные вектор-строки, но, имея в качестве основного пространство вектор-столбцов \mathbb{C}^n , нам будет удобно записать условие на левые собственные векторы в виде

$$y^*A = \mu y^*,$$

для $y \in \mathbb{C}^n$ и некоторого $\mu \in \mathbb{C}$. Применяя к этому соотношению эрмитово сопряжение, получим

$$A^*y = \bar{\mu}y,$$

т.е. левые собственные векторы матрицы A являются правыми собственными векторами эрмитово сопряжённой матрицы A^* . Эта простая взаимосвязь объясняет редкость самостоятельного использования

понятий левого и правого собственных векторов. Ясно, что при этом $\det(A^* - \bar{\mu}I) = 0$.

Исследуем подробнее так называемую *сопряжённую задачу* на собственные значения. Этим термином называют задачу нахождения собственных чисел и собственных векторов для эрмитово сопряжённой матрицы A^* :

$$A^*y = \varkappa y,$$

где $\varkappa \in \mathbb{C}$ — собственное значение матрицы A^* и $y \in \mathbb{C}^n$ — соответствующий собственный вектор. Как связаны между собой собственные значения и собственные векторы исходной A и сопряжённой A^* матриц? Для ответа на этот вопрос нам понадобится

Определение 3.2.1 *Два набора из одинакового количества векторов $\{r_1, r_2, \dots, r_m\}$ и $\{s_1, s_2, \dots, s_m\}$ в евклидовом или унитарном пространстве называются биортогональными, если $\langle r_i, s_j \rangle = 0$ при $i \neq j$.*

Приставка «би» в термине «биортогональность» означает, что введённое свойство относится к *двум* наборам векторов.

Ясно, что выполнение свойства биортогональности существенно зависит от порядка нумерации векторов в пределах каждого из наборов, так что в определении биортогональности неявно предполагается, что необходимые нумерации существуют и рассматриваемые наборы упорядочены в соответствии с ними. Нетрудно также понять, что если какой-либо набор векторов биортогонален сам себе, то он ортогонален в обычном смысле.

Предложение 3.2.2 *Собственные значения эрмитово-сопряжённых матриц попарно комплексно сопряжены друг другу. Собственные векторы эрмитово сопряжённых матриц биортогональны.*

Доказательство. Определитель матрицы, как известно, не меняется при её транспонировании, т. е. $\det A^T = \det A$. Комплексное сопряжение элементов матрицы влечёт комплексное сопряжение её определителя, $\det \bar{A} = \overline{\det A}$. Следовательно,

$$\begin{aligned} \det(A - \lambda I) &= \det(A - \lambda I)^T = \det(A^T - \lambda I) \\ &= \overline{\det(\overline{A^T - \lambda I})} = \overline{\det(A^* - \bar{\lambda}I)}. \end{aligned}$$

Отсюда мы можем заключить, что комплексное число z является корнем характеристического уравнения $\det(A - \lambda I) = 0$ матрицы A тогда и только тогда, когда ему сопряжённое \bar{z} является корнем уравнения $\det(A^* - \lambda I) = 0$, который является характеристическим для матрицы A^* . Это доказывает первое утверждение.

Пусть x и y — собственные векторы матриц A и A^* соответственно, а λ и \varkappa — отвечающие им собственные числа этих матриц. Для доказательства второго утверждения выпишем следующую цепочку преобразований:

$$\langle x, y \rangle = \frac{1}{\lambda} \langle \lambda x, y \rangle = \frac{1}{\lambda} \langle Ax, y \rangle = \frac{1}{\lambda} \langle x, A^* y \rangle = \frac{1}{\lambda} \langle x, \varkappa y \rangle = \frac{\bar{\varkappa}}{\lambda} \langle x, y \rangle.$$

Поэтому

$$\langle x, y \rangle - \frac{\bar{\varkappa}}{\lambda} \langle x, y \rangle = 0,$$

то есть

$$\langle x, y \rangle \left(1 - \frac{\bar{\varkappa}}{\lambda} \right) = 0.$$

Если x и y являются собственными векторами матриц A и A^* , отвечающими собственным значениям λ и \varkappa , которые не сопряжены комплексно друг другу, то в левой части полученного равенства второй сомножитель $(1 - \bar{\varkappa}/\lambda) \neq 0$. По этой причине необходимо $\langle x, y \rangle = 0$, что и требовалось доказать. ■

Обращаясь к определению правых и левых собственных векторов матрицы, можем утверждать, что если λ — правое собственное значение матрицы A , а μ — левое собственное значение, то $\bar{\lambda} = \bar{\mu}$. Иными словами, правые и левые собственные значения матрицы совпадают друг с другом, и потому их можно не различать. Что касается правых и левых собственных векторов матрицы, то они биортогональны друг другу.

Предложение 3.2.3 *Если λ — собственное число квадратной неособенной матрицы, то λ^{-1} — это собственное число обратной матрицы, отвечающее тому же собственному вектору.*

Доказательство. Если C — неособенная $n \times n$ -матрица и $Cv = \lambda v$, то $v = \lambda C^{-1}v$. Далее, так как $\lambda \neq 0$ в силу неособенности C , получаем отсюда $C^{-1}v = \lambda^{-1}v$. ■

3.2в Разложения матриц, использующие спектр

Квадратную матрицу вида

$$\begin{pmatrix} \alpha & 1 & & 0 \\ & \alpha & 1 & \\ & & \ddots & \ddots \\ 0 & & & \alpha & 1 \\ & & & & \alpha \end{pmatrix},$$

у которой по диагонали стоит α , на первой верхней побочной диагонали все единицы, а остальные элементы — нули, называют, как известно, *жордановой клеткой*, отвечающей значению α . Ясно, что α является собственным значением такой матрицы.

В линейной алгебре показывается, что с помощью подходящего преобразования подобия любая квадратная матрица может быть приведена к *жордановой канонической форме* — блочно-диагональной матрице, на главной диагонали которой стоят жордановы клетки, отвечающие собственным значениям рассматриваемой матрицы (см., к примеру, [7, 9, 23, 26, 38, 40, 50]). Иными словами для любой квадратной матрицы A существует такая неособенная матрица S , что

$$S^{-1}AS = J,$$

где

$$J = \left(\begin{array}{ccc|ccc} \lambda_1 & 1 & & & & \\ & \lambda_1 & \ddots & & 0 & 0 \\ & & \ddots & 1 & & \\ & & & \lambda_1 & & \\ \hline & 0 & & & \lambda_2 & 1 & \\ & & & & & \ddots & \ddots & \\ & & & & & & \lambda_2 & \\ \hline & 0 & & & & & & \ddots & \ddots & \end{array} \right),$$

а $\lambda_1, \lambda_2, \dots$ — собственные значения матрицы A .

Неприятной особенностью жордановой канонической формы является то, что она не зависит непрерывно от элементов матрицы, несмотря на то, что сами собственные значения матрицы непрерывно зависят от её элементов. Размеры жордановых клеток-блоков и их расположение вдоль диагонали могут скачкообразно меняться при изменении элементов матрицы. Это делает жорданову форму малоприменимой при решении многих практических задач, где входные данные носят приближённый и неточный характер.

Другое популярное разложение матриц, использующее информацию о спектре матрицы — это разложение Шура.

Пусть A — комплексная $n \times n$ -матрица и зафиксирован некоторый порядок её собственных значений $\lambda_1, \lambda_2, \dots, \lambda_n$. Существует такая унитарная $n \times n$ -матрица U , что матрица $T = U^*AU$ является верхней треугольной матрицей с диагональными элементами $\lambda_1, \lambda_2, \dots, \lambda_n$. Иными словами, любая квадратная матрица A унитарно эквивалентна треугольной матрице, в которой диагональные элементы являются собственными значениями для A , записанными в произвольном заранее заданном порядке. Если же A — это вещественная матрица и все её собственные значения вещественны, то U можно выбрать вещественной ортогональной матрицей. Представление

$$A = UTU^*$$

с верхней треугольной матрицей T и унитарными (ортогональными) матрицами U и U^* называют *разложением Шура* матрицы A . Оно в отличие от жордановой нормальной формы устойчиво к возмущениям элементов матрицы.

Для симметричных (эрмитовых в комплексном случае) матриц в выписанном представлении матрица T также должна быть симметричной (эрмитовой). Как следствие, в этом случае справедлив более сильный результат: с помощью ортогонального преобразования подобия любая матрица может быть приведена к диагональному виду, с собственными значениями по диагонали. Часто это представление называют спектральным разложением линейного оператора. Более общо, спектральное разложение — представление линейного оператора в виде линейной комбинации операторов проектирования на взаимно ортогональные оси.

3.2г Сингулярные числа и сингулярные векторы матрицы

Из результатов раздела 3.2б следует, что для определения собственных значений матрицы A и её левых и правых собственных векторов необходимо решить систему уравнений

$$\begin{cases} Ax = \lambda x, \\ y^* A = \lambda y^*. \end{cases} \quad (3.5)$$

Система уравнений (3.5) является «распавшейся»: в ней первые n уравнений и последние n уравнений не зависят друг от друга. Поэтому решать её также можно по частям, отдельно для x и отдельно для y , что обычно и делают на практике. Отметим, что для вещественных собственных чисел, когда $\lambda = \bar{\lambda}$, системе (3.5) после эрмитова сопряжения второй части можно придать следующий элегантный матричный вид

$$\begin{pmatrix} A & 0 \\ 0 & A^* \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \lambda \begin{pmatrix} x \\ y \end{pmatrix}. \quad (3.6)$$

Изменим соотношения (3.5), чтобы они «завязались» друг на друга, поменяв в правых частях векторы x и y :

$$\begin{cases} Ax = \sigma y, \\ y^* A = \sigma x^*. \end{cases} \quad (3.7)$$

Фигурально можно сказать, что при этом векторы x и y становятся «право-левыми» или «лево-правыми» собственными векторами матрицы A . Как мы увидим вскоре, аналоги собственных чисел матрицы, которые мы переобозначили через σ , также получают новое содержание.

Определение 3.2.2 *Неотрицательные вещественные скаляры σ , которые являются решениями системы матричных уравнений (3.7), называются сингулярными числами матрицы A . Удовлетворяющие системе (3.7) векторы x называются правыми сингулярными векторами матрицы A , а векторы y — левыми сингулярными векторами матрицы A .*

Отметим, что и система (3.7), и это определение имеют смысл уже для произвольных прямоугольных матриц, а не только для квадратных, как было в случае собственных значений и собственных векторов. Для $m \times n$ -матрицы A правые сингулярные векторы имеют размерность n , а левые — размерность m .

Если σ вещественно, то оно совпадает со своим комплексно-сопряжённым значением, $\bar{\sigma} = \sigma$, и потому, беря эрмитово сопряжение от второго уравнения (3.7), можем переписать систему уравнений для определения сингулярных чисел и сингулярных векторов в следующем виде:

$$\begin{cases} Ax = \sigma y, \\ A^*y = \sigma x. \end{cases} \quad (3.8)$$

Полезна также матричная форма системы (3.7):

$$\begin{pmatrix} 0 & A^* \\ A & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \sigma \begin{pmatrix} x \\ y \end{pmatrix}, \quad (3.9)$$

которая находится в красивой двойственности с системой (3.6). Если A — вещественная матрица, то векторы x и y также могут быть взяты вещественными, а система уравнений (3.9) для определения сингулярных чисел и векторов принимает ещё более простой вид:

$$\begin{pmatrix} 0 & A^\top \\ A & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \sigma \begin{pmatrix} x \\ y \end{pmatrix}.$$

Из уравнений (3.8)–(3.9) видно, что, в отличие от собственных значений, сингулярные числа характеризуют совместно как саму матрицу, так и её эрмитово-сопряжённую (транспонированную в вещественном случае).

Наша ближайшая цель — показать корректность Определения 3.2.2, то есть существование решений σ , x , y для системы уравнений (3.8)–(3.9) и наличие среди них неотрицательных σ .

Предложение 3.2.4 *Сингулярные числа матрицы A суть неотрицательные квадратные корни из собственных чисел матрицы A^*A или матрицы AA^* .*

Формулировка этого утверждения требует разъяснений, так как в случае прямоугольной $m \times n$ -матрицы A размеры квадратных матриц

A^*A и AA^* различны: первая из них — это $n \times n$ -матрица, а вторая $m \times m$ -матрица. Соответственно, количество собственных чисел у них будет различным.

Известно, что ранг произведения матриц не превосходит наименьшего из рангов перемножаемых матриц (см. [9, 23, 50]). Отсюда следует, что если $m < n$, то $n \times n$ -матрица A^*A имеет неполный ранг, не превосходящий m , а потому её собственные числа с $(m + 1)$ -го по n -ое — заведомо нулевые. Аналогично, если $m > n$, то неполный ранг, который не превосходит n , имеет $m \times m$ -матрица AA^* , и её собственные числа с $(n + 1)$ -го по m -ое равны нулю. Таким образом, для $m \times n$ -матрицы содержательный смысл имеет рассмотрение лишь $\min\{m, n\}$ штук сингулярных чисел, что устраняет вышеотмеченную кажущуюся неоднозначность.

Другой неочевидный момент формулировки Предложения 3.2.4 — взаимоотношение собственных чисел матриц A^*A и AA^* . Здесь можно вспомнить доказанный выше общий результат линейной алгебры — Предложение 3.2.1, — о совпадении ненулевых точек спектра произведений двух матриц, взятых в различном порядке. Впрочем, для частного случая матриц A^*A и AA^* этот момент будет обоснован в следующем ниже доказательстве.

Доказательство. Умножая обе части второго уравнения из (3.8) на σ , получим $A^*(\sigma y) = \sigma^2 x$. Затем подставим сюда значение σy из первого уравнения (3.8): $A^*Ax = \sigma^2 x$.

С другой стороны, умножая на σ обе части первого уравнения (3.8), получим $A(\sigma x) = \sigma^2 y$. Подставив сюда значение σx из второго уравнения (3.7), получим $AA^*y = \sigma^2 y$. Иными словами, числа σ^2 являются собственными числами как для A^*A , так и для AA^* .

Покажем теперь, что собственные значения у матриц A^*A и AA^* неотрицательны, чтобы иметь возможность извлекать из них квадратные корни для окончательного определения σ . Очевидно, это достаточно сделать лишь для одной из выписанных матриц, так как для другой рассуждения совершенно аналогичны.

Коль скоро матрица A^*A эрмитова, любое её собственное значение λ вещественно. Кроме того, если u — соответствующий собственный вектор, то $0 \leq (Au)^*(Au) = u^*(A^*Au) = u^*\lambda u = \lambda u^*u$, откуда в силу $u^*u > 0$ следует $\lambda \geq 0$.

Для завершения доказательства осталось продемонстрировать, что арифметические квадратные корни из собственных значений матриц

A^*A и AA^* вместе с их собственными векторами удовлетворяют системе уравнений (3.8)–(3.9).

Пусть u — собственный вектор матрицы A^*A , отвечающий собственному числу λ , так что

$$A^*Au = \lambda u,$$

причём $\lambda \geq 0$ в силу ранее доказанного. Обозначим $y := Au$ и $x := \sqrt{\lambda}u$. Тогда $\lambda u = \sqrt{\lambda}x$, и потому

$$\begin{aligned} Ax &= A(\sqrt{\lambda}u) = \sqrt{\lambda}Au = \sqrt{\lambda}y, \\ A^*y &= A^*Au = \lambda u = \sqrt{\lambda}x, \end{aligned}$$

так что система (3.7)–(3.9) удовлетворяется при $\sigma = \sqrt{\lambda}$ с выбранными векторами x и y .

Аналогично, если v — собственный вектор матрицы AA^* , отвечающий её собственному числу μ , то

$$AA^*v = \mu v,$$

причём $\mu \geq 0$. Обозначим $x := A^*v$ и $y := \sqrt{\mu}v$. Тогда $\mu v = \sqrt{\mu}y$, и потому

$$\begin{aligned} Ax &= AA^*v = \mu v = \sqrt{\mu}y, \\ A^*y &= A^*(\sqrt{\mu}v) = \sqrt{\mu}A^*v = \sqrt{\mu}x, \end{aligned}$$

так что система (3.8)–(3.9) действительно удовлетворяется при $\sigma = \sqrt{\mu}$ с выбранными векторами x и y . ■

Итак, задаваемые Определением 3.2.2 сингулярные числа вещественной или комплексной $m \times n$ -матрицы — это набор из $\min\{m, n\}$ неотрицательных вещественных чисел, которые обычно нумеруют в порядке убывания:

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m, n\}} \geq 0.$$

Таким образом, $\sigma_1 = \sigma_1(A)$ — это наибольшее сингулярное число матрицы A . Мы будем также обозначать наибольшее и наименьшее сингулярные числа матрицы посредством $\sigma_{\max}(A)$ и $\sigma_{\min}(A)$.

Из доказательства Предложения 3.2.4 следует также, что правыми сингулярными векторами матрицы A являются правые собственные векторы матрицы A^*A , а левыми сингулярными векторами матрицы A

— левые собственные векторы для A^*A или, что равносильно, эрмитово сопряжённые правых собственных векторов матрицы AA^* . Отметим также, что как левые, так и правые сингулярные векторы суть ортогональные системы векторов, коль скоро они являются собственными векторами эрмитовых матриц A^*A и AA^* .

Пример 3.2.1 Пусть A — это 1×1 -матрица, т.е. просто некоторое число a , вещественное или комплексное. Ясно, что единственное собственное число такой матрицы равно самому a . Сингулярное число у A также всего одно, и оно равно $\sqrt{a^*a} = |a|$.

Пусть $A = (a_1, a_2, \dots, a_n)^T$ — это $n \times 1$ -матрица, т.е. просто вектор-столбец. Тогда матрица $A^T A$ является числом $a_1^2 + a_2^2 + \dots + a_n^2$, и поэтому единственное сингулярное число матрицы A равно евклидовой норме вектора $(a_1, a_2, \dots, a_n)^T$. То же самое верно для $1 \times n$ -матрицы, то есть вектор-строки (a_1, a_2, \dots, a_n) . ■

Пример 3.2.2 Для единичной матрицы I все сингулярные числа очевидно равны единицам.

Но все единичные сингулярные числа имеет не только единичная матрица. Если Q — унитарная комплексная матрица (ортогональная в вещественном случае), то $Q^*Q = I$, и потому все сингулярные числа для Q также равны единицам. ■

Пример 3.2.3 Для 2×2 -матрицы

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \quad (3.10)$$

нетрудно выписать характеристическое уравнение

$$\det \begin{pmatrix} 1 - \lambda & 2 \\ 3 & 4 - \lambda \end{pmatrix} = \lambda^2 - 4\lambda - 2 = 0,$$

и найти его корни $\frac{1}{2}(5 \pm \sqrt{33})$ — собственные значения матрицы, приблизительно равные -0.372 и 5.372 . Для определения сингулярных чисел образуем

$$A^T A = \begin{pmatrix} 10 & 14 \\ 14 & 20 \end{pmatrix},$$

и вычислим её собственные значения. Они равны $15 \pm \sqrt{221}$, и потому

получается, что сингулярные числа матрицы A суть $\sqrt{15 \pm \sqrt{221}}$, т. е. примерно 0.366 и 5.465 (с точностью до трёх знаков после запятой).

С другой стороны, для матрицы

$$\begin{pmatrix} 1 & 2 \\ -3 & 4 \end{pmatrix}, \quad (3.11)$$

которая отличается от матрицы (3.10) лишь противоположным знаком элемента на месте $(2, 1)$, собственные значения — это комплексно-сопряжённая пара $\frac{1}{2}(5 \pm i\sqrt{15}) \approx 2.5 \pm 1.936i$, а сингулярные числа суть $\sqrt{15 \pm \sqrt{125}}$, т. е. приблизительно 1.954 и 5.117. ■

Можно заметить, что максимальные сингулярные числа рассмотренных матриц превосходят наибольшие из модулей их собственных чисел. Мы увидим ниже (см. §3.3ж), что это не случайно, и наибольшее сингулярное число всегда не меньше, чем максимум модулей собственных чисел матрицы.

Рассмотрим вопрос о том, как связаны сингулярные числа для взаимно обратных матриц.

Предложение 3.2.5 *Если σ — сингулярное число неособенной квадратной матрицы, то σ^{-1} — это сингулярное число обратной матрицы.*

Доказательство. Вспомним, что собственные числа взаимно обратных матриц обратны друг другу. Применяя это соображение к матрице A^*A , можем заключить, что если $\lambda_1, \lambda_2, \dots, \lambda_n$ — её собственные значения, то у обратной матрицы $(A^*A)^{-1} = A^{-1}(A^*)^{-1}$ собственными значениями являются $\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_n^{-1}$. Но $A^{-1}(A^*)^{-1} = A^{-1}(A^{-1})^*$, а потому в силу Предложения 3.2.4 выписанные числа $\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_n^{-1}$ образуют набор квадратов сингулярных чисел матрицы A^{-1} . Это и требовалось показать. ■

3.2д Сингулярное разложение матриц

Важнейший результат, касающийся сингулярных чисел и сингулярных векторов матриц, который служит одной из основ их широкого применения в разнообразных вопросах математического моделирования — это

Теорема 3.2.1 (теорема о сингулярном разложении матрицы)

Для любой комплексной $m \times n$ -матрицы A существуют унитарные $m \times m$ -матрица U и $n \times n$ -матрица V , такие что

$$A = U \Sigma V^* \quad (3.12)$$

с диагональной $m \times n$ -матрицей

$$\Sigma = \begin{pmatrix} \sigma_1 & 0 & 0 & \cdots & 0 \\ 0 & \sigma_2 & 0 & \cdots & 0 \\ 0 & 0 & \sigma_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \end{pmatrix},$$

где $\sigma_1, \sigma_2, \dots, \sigma_{\min\{m,n\}}$ — сингулярные числа матрицы A , а столбцы матриц U и V являются соответственно левыми и правыми сингулярными векторами матрицы A .

Представление (3.12) называется *сингулярным разложением матрицы* A . Если A — вещественная матрица, то U и V являются также вещественными ортогональными матрицами, и сингулярное разложение принимает вид

$$A = U \Sigma V^T.$$

Для квадратных матриц доказательство сингулярного разложения может быть легко выведено из известного полярного разложения матрицы, т. е. её представления в виде $A = QS$, где в комплексном случае Q — унитарная матрица, S — эрмитова, а в вещественном случае Q — ортогональная, S — симметричная (см., к примеру, [9, 23, 50]). Рассмотрим подробно общий комплексный случай.

Как известно, любую эрмитову матрицу можно унитарными преобразованиями подобия привести к диагональному виду, так что $S = T^*DT$, где T — унитарная, а D — диагональная. Поэтому $A = (QT^*)DT$. Это уже почти требуемое представление для A , поскольку произведение унитарных матриц Q и T^* тоже унитарно. Нужно лишь убедиться в том, что по диагонали в D стоят сингулярные числа матрицы A .

Исследуем произведение A^*A :

$$\begin{aligned} A^*A &= ((QT^*)DT)^*((QT^*)DT) \\ &= T^*D^*(QT^*)^*(QT^*)DT \\ &= T^*D^*DT = T^*D^2T. \end{aligned}$$

Как видим, матрица A^*A подобна диагональной матрице D^2 , их собственные числа поэтому совпадают. Следовательно, собственные числа A^*A суть квадраты диагональных элементов D . Это и требовалось доказать.

В общем случае доказательство Теоремы 3.2.1 не очень сложно и может быть найдено, к примеру, в книгах [11, 38, 40]. Фактически, этот результат показывает, как с помощью сингулярных чисел матрицы элегантно представляется действие соответствующего линейного оператора из одного векторного пространства в другое. Именно, для любого линейного отображения можно выбрать ортонормированный базис в пространстве области определения и ортонормированный базис в пространстве области значений так, чтобы в этих базисах рассматриваемое отображение представлялось растяжением вдоль координатных осей. Сингулярные числа матрицы оказываются, как правило, адекватным инструментом её исследования, когда соответствующее линейное отображение действует из одного векторного пространства в другое, возможно с отличающимися друг от друга размерностями. Собственные числа матрицы полезны при изучении линейного преобразования векторного пространства в пространство той же размерности, в частности, самого в себя.

Другие примеры применения сингулярных чисел и сингулярных векторов матрицы рассматриваются ниже в §3.4.

Сингулярное разложение матриц впервые возникло в конце XIX века в трудах Э. Бельтрами и К. Жордана, но термин *valeurs singulières* — «сингулярные значения» — впервые был использован французским математиком Э. Пикаром около 1910 года в работе по интегральным уравнениям (см. [92]). Задача нахождения сингулярных чисел и сингулярных векторов матриц, последняя из списка на стр. 200, по видимости является частным случаем третьей задачи, относящейся к нахождению собственных чисел и собственных векторов. Но вычисление сингулярных чисел и векторов матриц сделалось в настоящее время очень важным как в теории, так и в приложениях вычислительной линейной

алгебры. С другой стороны, соответствующие численные методы весьма специализированы, так что эта задача в общем списке задач уже выделяется отдельным пунктом.

Комментируя современный список задач вычислительной линейной алгебры из §3.1, можно также отметить, что на первые места в нём выдвинулась линейная задача о наименьших квадратах. А некоторые старые и популярные ранее задачи как бы отошли на второй план, что стало отражением значительных изменений в вычислительных технологиях решения задач математического моделирования. Это естественный процесс, в котором большую роль сыграло развитие современной вычислительной техники. Следует быть готовым к подобным изменениям и в будущем.

3.2e Матрицы с диагональным преобладанием

В приложениях линейной алгебры и теории матриц часто возникают матрицы, в которых диагональные элементы в том или ином смысле «преобладают» над остальной, недиагональной частью матрицы. Это обстоятельство может быть, к примеру, следствием особенностей рассматриваемой математической модели, в которой связи составляющих её частей с самими собой (они и выражаются диагональными элементами) сильнее, чем с остальными. Такие матрицы обладают рядом замечательных свойств, изложению одного из которых и посвящён этот пункт. Следует отметить, что сам смысл, вкладываемый в понятие «преобладания» может быть различен, и ниже мы рассмотрим простейший и наиболее популярный.

Определение 3.2.3 *Квадратную $n \times n$ -матрицу $A = (a_{ij})$ называют матрицей с диагональным преобладанием, если для любого $i = 1, 2, \dots, n$ имеет место*

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|. \quad (3.13)$$

Матрицы, удовлетворяющие этому определению, некоторые авторы называют матрицами со «строгим диагональным преобладанием». Со своей стороны, мы будем говорить, что матрица $A = (a_{ij})$ имеет *нестрогое диагональное преобладание* в случае выполнения неравенств

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}| \quad (3.14)$$

для любого $i = 1, 2, \dots, n$. Иногда в связи с условиями (3.13) и (3.14) необходимо уточнять, что речь идёт о диагональном преобладании «по строкам», поскольку имеет также смысл диагональное преобладание «по столбцам», которое определяется совершенно аналогичным образом.

Теорема 3.2.2 (признак Адамара) *Матрица с диагональным преобладанием неособенна.*

Доказательство. Предположим, что, вопреки доказываемому, рассматриваемая матрица $A = (a_{ij})$ является особенной. Тогда для некоторого ненулевого n -вектора $y = (y_1, y_2, \dots, y_n)^T$ выполняется равенство $Ay = 0$, т. е.

$$\sum_{j=1}^n a_{ij}y_j = 0, \quad i = 1, 2, \dots, n. \quad (3.15)$$

Выберем среди компонент вектора y ту, которая имеет наибольшее абсолютное значение. Пусть она имеет номер ν , так что $|y_\nu| = \max_{1 \leq j \leq n} |y_j|$, причём $|y_\nu| > 0$ в силу сделанного выше предположения о том, что $y \neq 0$. Следствием ν -го из равенств (3.15) является соотношение

$$-a_{\nu\nu}y_\nu = \sum_{j \neq \nu} a_{\nu j}y_j,$$

которое влечёт цепочку оценок

$$\begin{aligned} |a_{\nu\nu}| |y_\nu| &= \left| \sum_{j \neq \nu} a_{\nu j}y_j \right| \leq \sum_{j \neq \nu} |a_{\nu j}| |y_j| \\ &\leq \left(\max_{1 \leq j \leq n} |y_j| \right) \sum_{j \neq \nu} |a_{\nu j}| = |y_\nu| \sum_{j \neq \nu} |a_{\nu j}|. \end{aligned}$$

Сокращая теперь обе части полученного неравенства на $|y_\nu| > 0$, будем иметь

$$|a_{\nu\nu}| \leq \sum_{j \neq \nu} |a_{\nu j}|,$$

что противоречит неравенствам (3.13), т. е. наличию, по условию теоремы, диагонального преобладания в матрице A . Итак, A действительно должна быть неособенной матрицей. ■

Доказанный выше результат часто именуют также «теоремой Леви-Деспланка» (см., к примеру, [41, 50]), но мы придерживаемся здесь терминологии, принятой в [9, 79]. В книге М. Пароди [79] можно прочитать, в частности, некоторые сведения об истории вопроса.

Внимательное изучение доказательства признака Адамара показывает, что в нём нигде не использовался факт принадлежности элементов матрицы и векторов какому-то конкретному числовому полю — \mathbb{R} или \mathbb{C} . Таким образом, признак Адамара справедлив и для комплексных матриц. Кроме того, он может быть отчасти обобщен на матрицы, удовлетворяющие нестрогому диагональному преобладанию (3.14).

Вещественная или комплексная $n \times n$ -матрица $A = (a_{ij})$ называется *разложимой*, если существует разбиение множества $\{1, 2, \dots, n\}$ первых n натуральных чисел на два непересекающихся подмножества I и J , таких что $a_{ij} = 0$ при $i \in I$ и $j \in J$. Эквивалентное определение: матрица $A \in \mathbb{R}^{n \times n}$ разложима, если путём перестановок строк и столбцов она может быть приведена к блочно-треугольному виду

$$\begin{pmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{pmatrix}$$

с квадратными блоками A_{11} и A_{22} . Матрицы, не являющиеся разложимыми, называются *неразложимыми*. Важнейший пример неразложимых матриц — это матрицы, все элементы которых не равны нулю, в частности, положительны.

Обобщением признака Адамара является

Теорема 3.2.3 (теорема Гаусски) *Если для квадратной неразложимой матрицы A выполнены условия нестрогого диагонального преобладания (3.14), причём хотя бы одно из этих неравенств выполнено строго, то матрица A неособенна.*

Доказательство можно найти, к примеру, в [9].

3.3 Нормы векторов и матриц

3.3а Векторные нормы

Норму можно рассматривать как обобщение на многомерный и абстрактный случаи понятия абсолютной величины числа. Вообще, и

норма, и абсолютная величина являются понятиями, которые формализуют интуитивно ясное свойство «размера» объекта, его «величины», т. е. того, насколько он мал или велик безотносительно к его расположению в пространстве или к другим второстепенным качествам. Такова, например, длина вектора как направленного отрезка в привычном нам евклидовом пространстве.

Формальное определение даётся следующим образом:

Определение 3.3.1 *Нормой в вещественном или комплексном линейном векторном пространстве X называется вещественнозначная функция $\|\cdot\|$, удовлетворяющая следующим свойствам (называемым аксиомами нормы):*

$$(ВН1) \quad \|a\| \geq 0 \quad \text{для любого } a \in X, \text{ причём } \|a\| = 0 \Leftrightarrow a = 0$$

— неотрицательность,

$$(ВН2) \quad \|\alpha a\| = |\alpha| \cdot \|a\| \quad \text{для любых } a \in X \text{ и } \alpha \in \mathbb{R} \text{ или } \mathbb{C}$$

— абсолютная однородность,

$$(ВН3) \quad \|a + b\| \leq \|a\| + \|b\| \quad \text{для любых } a, b, c \in X$$

— «неравенство треугольника».

Само пространство X с нормой называется при этом нормированным линейным пространством.

Далее в качестве конкретных линейных векторных пространств у нас, как правило, всюду рассматриваются пространства \mathbb{R}^n или \mathbb{C}^n .

Не все нормы, удовлетворяющие выписанным аксиомам одинаково практичны, и часто от нормы требуют выполнения ещё тех или иных дополнительных условий. К примеру, удобно иметь дело с *абсолютной нормой*, значение которой зависит лишь от абсолютных значений компонент векторов. В общем случае норма вектора этому условию может и не удовлетворять.

Приведём примеры наиболее часто используемых норм векторов в \mathbb{R}^n и \mathbb{C}^n . Если $a = (a_1, a_2, \dots, a_n)^\top$, то обозначим

$$\|a\|_1 := \sum_{i=1}^n |a_i|,$$

$$\|a\|_2 := \left(\sum_{i=1}^n |a_i|^2 \right)^{1/2},$$

$$\|a\|_\infty := \max_{1 \leq i \leq n} |a_i|.$$

Вторая из этих норм часто называется *евклидовой*, а третья — *чебышёвской* или *максимум-нормой*. Евклидова норма вектора, как направленного отрезка, — это его обычная длина, в связи с чем евклидову норму часто называют также *длиной вектора*. Нередко можно встретить и другие названия рассмотренных норм.

Замечательность евклидовой нормы $\|\cdot\|_2$ состоит в том, что она порождается стандартным скалярным произведением $\langle \cdot, \cdot \rangle$ в \mathbb{R}^n или \mathbb{C}^n . Более точно, если скалярное произведение задаётся как (3.1) или (3.2), то $\|a\|_2 = \sqrt{\langle a, a \rangle}$. Иными словами, 2-норма является составной частью более богатой и содержательной структуры на пространствах \mathbb{R}^n и \mathbb{C}^n , чем мы будем неоднократно пользоваться. Напомним также *неравенство Коши-Буняковского*

$$|\langle a, b \rangle| \leq \|a\|_2 \|b\|_2. \quad (3.16)$$

Нормы $\|\cdot\|_1$ и $\|\cdot\|_2$ — это частные случаи более общей конструкции *p-нормы*

$$\|a\|_p = \left(\sum_{i=1}^n |a_i|^p \right)^{1/p} \quad \text{для } p \geq 1,$$

которую называют также *гёльдеровой нормой* (по имени О.Л. Гёльдера). Неравенство треугольника для неё имеет вид

$$\left(\sum_{i=1}^n |a_i + b_i|^p \right)^{1/p} \leq \left(\sum_{i=1}^n |a_i|^p \right)^{1/p} + \left(\sum_{i=1}^n |b_i|^p \right)^{1/p},$$

оно называется *неравенством Минковского* и имеет самостоятельное значение в различных разделах математики. Чебышёвская норма также может быть получена из *p-нормы* с помощью предельного перехода по $p \rightarrow \infty$, что и объясняет индекс « ∞ » в её обозначении.

В самом деле,

$$\left(\sum_{i=1}^n |a_i|^p \right)^{1/p} \leq \left(n \left(\max_{1 \leq i \leq n} |a_i| \right)^p \right)^{1/p} = n^{1/p} \max_{1 \leq i \leq n} |a_i|.$$

С другой стороны,

$$\left(\sum_{i=1}^n |a_i|^p \right)^{1/p} \geq \left(\left(\max_{1 \leq i \leq n} |a_i| \right)^p \right)^{1/p} = \max_{1 \leq i \leq n} |a_i|,$$

так что в целом

$$\max_{1 \leq i \leq n} |a_i| \leq \left(\sum_{i=1}^n |a_i|^p \right)^{1/p} \leq n^{1/p} \max_{1 \leq i \leq n} |a_i|.$$

При переходе в этом двойном неравенстве к пределу по $p \rightarrow \infty$ оценки снизу и сверху сливаются, и потому действительно

$$\lim_{p \rightarrow \infty} \left(\sum_{i=1}^n |a_i|^p \right)^{1/p} = \max_{1 \leq i \leq n} |a_i|.$$

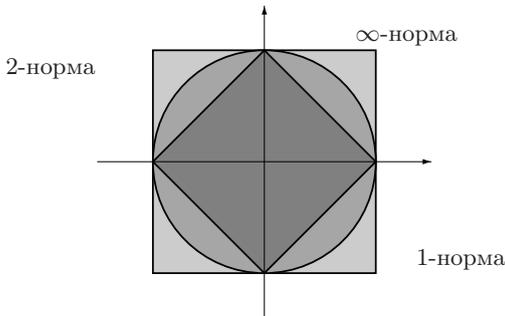


Рис. 3.5. Шары единичного радиуса в различных нормах.

В нормированном пространстве \mathcal{X} шаром радиуса r с центром в точке a называется множество $\{x \in \mathcal{X} \mid \|x - a\| \leq r\}$. Геометрически наглядное представление о норме даётся её единичным шаром, т. е. множеством $\{x \mid \|x\| \leq 1\}$. На Рис. 3.5 нарисованы единичные шары

для рассмотренных выше норм в \mathbb{R}^2 . Из аксиом нормы вытекает, что единичный шар любой нормы — это множество в линейном векторном пространстве, которое выпукло (следствие неравенства треугольника) и *уравновешено*, т. е. инвариантно относительно умножения на любой скаляр α с $|\alpha| \leq 1$ (следствие абсолютной однородности).

Нередко используются взвешенные (масштабированные) варианты норм векторов, в выражениях для которых каждая компонента берётся с каким-то положительным весовым коэффициентом, отражающим его индивидуальный вклад в рассматриваемую модель. В частности, взвешенная чебышёвская норма определяется для положительного весового вектора $(\gamma_1, \gamma_2, \dots, \gamma_n)$, $\gamma_i > 0$, как

$$\|a\|_{\infty, \gamma} = \max_{1 \leq i \leq n} |\gamma_i a_i|.$$

Её единичные шары — различные прямоугольные брусы с гранями, параллельными координатным осям, т. е. прямые произведения интервалов вещественной оси (см. Рис. 3.6). Они являются частным случаем многомерных интервалов, и в связи с этим обстоятельством взвешенная чебышёвская норма популярна в интервальном анализе.

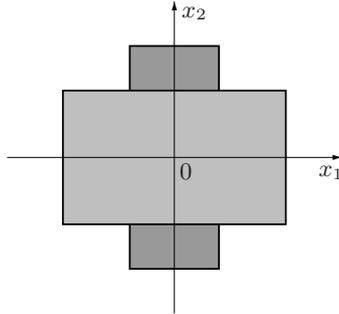


Рис. 3.6. Шары единичного радиуса во взвешенных чебышёвских нормах.

Обобщением конструкции взвешенных норм может служить норма, связанная с некоторой фиксированной неособенной матрицей. Именно, если $\|\cdot\|$ — какая-либо векторная норма в \mathbb{R}^n или \mathbb{C}^n , а S — неособенная $n \times n$ -матрица, то можно определить норму векторов как $\|x\|_S = \|Sx\|$. Нетрудно проверить, что все аксиомы векторной нормы удовлетворяются. Мы воспользуемся такой нормой ниже в §3.9б.

3.36 Топология на векторных пространствах

Говорят, что на множестве X задана *топологическая структура*, или просто *топология*⁵, если в X выделен класс подмножеств, содержащий вместе с каждым набором множеств их объединение, и вместе с каждым конечным набором множеств — их пересечение. Множество, снабжённое топологической структурой, называется *топологическим пространством*, а множества выделенного класса — *открытыми множествами*. Подмножество топологического пространства называется *замкнутым*, если его дополнение открыто.

Окрестностью точки в топологическом пространстве называется всякое открытое множество, содержащее эту точку. Окрестностью подмножества топологического пространства называется всякое открытое множество, содержащее это подмножество. Задание окрестностей точек и множеств позволяет определять близость одного элемента множества к другому, предельные переходы, сходимости и т. п. понятия. Топологическую структуру (топологию) можно задавать различными способами, например, простым описанием того, какие именно множества считаются открытыми.

В практике математического моделирования более распространено задание топологии не сформулированным выше абстрактным способом, а при помощи функции расстояния (метрики) или же с помощью различных норм. Преимущество этого пути состоит в том, что мы получаем в своё распоряжение количественную меру близости рассматриваемых объектов. При этом открытыми множествами считаются такие множества, каждая точка которых принадлежит множеству вместе с некоторым шаром с центром в этой точке.

Как известно, на нормированном пространстве X расстояние (метрика) между элементами a и b может быть естественно задано как

$$\text{dist}(a, b) = \|a - b\|,$$

т. е. как «величина различия» элементов a и b . Непосредственной проверкой легко убедиться, что для введённого таким образом расстояния dist выполняются все аксиомы расстояния (мы приводили их ранее на стр. 43). Таким образом, нормы будут нужны нам как сами по себе,

⁵Топологией называется также математическая дисциплина, изучающая, главным образом, свойства объектов, инвариантные относительно непрерывных отображений (см., к примеру, [60]). Ниже даётся очень краткий обзор основных идей топологии, предназначенный, скорее, для напоминания или увлечения читателя.

для оценивания «величины» тех или иных объектов, так и для измерения «отклонения» одного вектора от другого. Кроме того, задание нормы на некотором линейном векторном пространстве X автоматически определяет на нём и топологию, т. е. запас открытых и замкнутых множеств, структуру близости, с помощью которой можно будет, в частности, выполнять предельные переходы. Более точно, введём следующее

Определение 3.3.2 *Говорят, что в нормированном пространстве X с нормой $\|\cdot\|$ переменная $a \in X$ сходится к пределу a^* по норме (относительно рассматриваемой нормы), если $\|a - a^*\| \rightarrow 0$.*

Нормы в линейном векторном пространстве называются *топологически эквивалентными* (или просто *эквивалентными*), если эквивалентны порождаемые ими топологии, т. е. любое открытое (замкнутое) относительно одной нормы множество является открытым (замкнутым) также в другой норме, и наоборот. При условии эквивалентности норм, в частности, предельный переход в одной из них влечёт существование предела в другой, и наоборот. Из математического анализа известен простой критерий эквивалентности двух норм (см., к примеру, [52]):

Предложение 3.3.1 *Нормы $\|\cdot\|'$ и $\|\cdot\|''$ на линейном векторном пространстве X эквивалентны тогда и только тогда, когда существуют такие положительные константы C_1 и C_2 , что для любых $a \in X$*

$$C_1 \|a\|' \leq \|a\|'' \leq C_2 \|a\|'. \quad (3.17)$$

Формулировка этого предложения имеет кажущуюся асимметрию, так как для значений одной из эквивалентных норм предъявляется двусторонняя «вилка» из значений другой нормы с подходящими множителями-константами. Но нетрудно видеть, что из (3.17) немедленно следует

$$\frac{1}{C_2} \|a\|'' \leq \|a\|' \leq \frac{1}{C_1} \|a\|'',$$

так что существование «вилки» для одной нормы автоматически подразумевает существование аналогичной «вилки» и для другой. C_1 и C_2 обычно называют *константами эквивалентности* норм $\|\cdot\|'$ и $\|\cdot\|''$.

Содержательный смысл Предложения 3.3.1 совершенно прозрачен. Если $C_1 \|a\|' \leq \|a\|''$, то в любой шар ненулевого радиуса в норме $\|\cdot\|''$

можно вложить некоторый шар в норму $\|\cdot\|'$. Если же $\|a\|'' \leq C_2\|a\|'$, то верно и обратное: в любой шар ненулевого радиуса относительно нормы $\|\cdot\|'$ можно поместить какой-то шар ненулевого радиуса относительно нормы $\|\cdot\|''$. Как следствие, множество, открытое относительно одной нормы, будет также открытым относительно другой, и наоборот. По этой причине одинаковыми окажутся запасы окрестностей любой точки, так что топологические структуры, порождаемые этими двумя нормами, будут эквивалентны друг другу.

Предложение 3.3.2 *В векторных пространствах \mathbb{R}^n или \mathbb{C}^n*

$$\|a\|_2 \leq \|a\|_1 \leq \sqrt{n} \|a\|_2,$$

$$\|a\|_\infty \leq \|a\|_2 \leq \sqrt{n} \|a\|_\infty,$$

$$\frac{1}{n} \|a\|_1 \leq \|a\|_\infty \leq \|a\|_1,$$

т. е. векторные 1-норма, 2-норма и ∞ -норма эквивалентны друг другу.

Доказательство. Справедливость правого из первых неравенств следует из неравенства Коши-Буняковского (3.16), применённого к случаю $b = (\operatorname{sgn} a_1, \operatorname{sgn} a_2, \dots, \operatorname{sgn} a_n)^\top$. Для обоснования левого из первых неравенств заметим, что в силу определений 2-нормы и 1-нормы

$$\|a\|_2^2 = |a_1|^2 + |a_2|^2 + \dots + |a_n|^2,$$

$$\begin{aligned} \|a\|_1^2 &= |a_1|^2 + |a_2|^2 + \dots + |a_n|^2 \\ &\quad + 2|a_1a_2| + 2|a_1a_3| + \dots + 2|a_{n-1}a_n|, \end{aligned}$$

и все слагаемые $2|a_1a_2|, 2|a_1a_3|, \dots, 2|a_{n-1}a_n|$ неотрицательны. В частности, равенство $\|a\|_2^2 = \|a\|_1^2$ и ему равносильное $\|a\|_2 = \|a\|_1$ возможны лишь в случае, когда у вектора a все компоненты равны нулю за исключением одной.

Обоснование остальных неравенств дается следующими несложными

ми выкладками:

$$\begin{aligned} \|a\|_2 &= \sqrt{|a_1|^2 + |a_2|^2 + \dots + |a_n|^2} \\ &\geq \sqrt{\max_i |a_i|^2} = \max_i |a_i| = \|a\|_\infty, \\ \|a\|_2 &= \sqrt{|a_1|^2 + |a_2|^2 + \dots + |a_n|^2} \\ &\leq \sqrt{n \max_i |a_i|^2} = \sqrt{n} \max_i |a_i| = \sqrt{n} \|a\|_\infty, \\ \|a\|_\infty &= \max_i |a_i| \\ &\leq |a_1| + |a_2| + \dots + |a_n| = \|a\|_1, \\ \|a\|_1 &= |a_1| + |a_2| + \dots + |a_n| \\ &\leq n \max_i |a_i| \leq n \|a\|_\infty. \end{aligned}$$

Нетрудно видеть, что все эти неравенства достижимы (точные). ■

Доказанный выше вывод об эквивалентности конкретных норм является частным случаем общего результата математического анализа: *в конечномерном линейном векторном пространстве все нормы топологически эквивалентны друг другу* (см., к примеру, [20, 40, 50]). Но содержание Предложения 3.3.2 состоит ещё и в указании конкретных констант эквивалентности норм, от которых существенно зависят различные числовые оценки и вытекающие из них действия по решению тех или иных задач.

Любой вектор однозначно представляется своим разложением по какому-то фиксированному базису линейного пространства, или, иными словами, своими компонентами-числами в этом базисе. В связи с этим помимо определённой выше сходимости по норме имеет смысл рассматривать “покомпонентную сходимость”, при которой один вектор считается сходящимся к другому тогда и только тогда, когда все компоненты первого вектора сходятся к соответствующим компонентам второго. Формализацией этих соображений является

Определение 3.3.3 *Говорят, что переменная $a \in \mathcal{X}$ сходится к пределу a^* покомпонентно (покомпонентным образом) относительно некоторого базиса, если для каждого индекса i имеет место сходимость соответствующей компоненты $a_i \rightarrow a_i^*$ в \mathbb{R} или \mathbb{C} .*

Интересен вопрос о том, как соотносятся между собой сходимость по норме и сходимость всех компонент вектора.

Предложение 3.3.3 *В конечномерных линейных векторных пространствах сходимость по норме и покомпонентная сходимость векторов равносильны друг другу.*

Доказательство. Пусть a — n -мерная векторная переменная, которая сходится к пределу a^* в покомпонентном смысле относительно базиса $\{e_i\}_{i=1}^n$. Тогда, разлагая a и a^* в этом базисе, получаем

$$\begin{aligned} \|a - a^*\| &= \left\| \sum_{i=1}^n a_i e_i - \sum_{i=1}^n a_i^* e_i \right\| = \left\| \sum_{i=1}^n (a_i - a_i^*) e_i \right\| \\ &\leq \sum_{i=1}^n \|(a_i - a_i^*) e_i\| = \sum_{i=1}^n |a_i - a_i^*| \|e_i\|. \end{aligned}$$

Как следствие, если a_i сходятся к a_i^* для любого индекса $i = 1, 2, \dots, n$, то и $\|a - a^*\| \rightarrow 0$.

Обратно, пусть имеет место сходимость a к a^* по норме. Из факта эквивалентности различных норм следует существование такой положительной константы C , что

$$\max_i |a_i - a_i^*| = \|a - a^*\|_\infty \leq C \|a - a^*\|.$$

Поэтому при $\|a - a^*\| \rightarrow 0$ обязательно должна быть сходимость компонент a_i к a_i^* для всех индексов i . ■

Хотя сходимость по норме и покомпонентная сходимость равносильны друг другу, в различных ситуациях часто бывает удобнее воспользоваться какой-нибудь одной из них. Норма является одним числом, указывающим на степень близости к пределу, и работать с ней поэтому проще. Но рассмотрение сходимости в покомпонентном смысле позволяет расчленивать задачу на отдельные компоненты, что также нередко упрощает рассуждения.

Покажем непрерывность сложения и умножения на скаляр относительно нормы. Пусть $a \rightarrow a^*$ и $b \rightarrow b^*$, так что $\|a - a^*\| \rightarrow 0$ и

$\|b - b^*\| \rightarrow 0$. Тогда

$$\begin{aligned} \|(a + b) - (a^* + b^*)\| &= \|(a - a^*) + (b - b^*)\| \leq \|a - a^*\| + \|b - b^*\| \rightarrow 0, \\ \|\alpha a - \alpha a^*\| &= \|\alpha(a - a^*)\| = |\alpha| \|a - a^*\| \rightarrow 0 \end{aligned}$$

для любого скаляра α .

Умножение на матрицу также непрерывно в конечномерном линейном векторном пространстве. Если A — $m \times n$ -матрица и b — такой n -вектор, что $b \rightarrow b^*$, то, зафиксировав индекс $i \in \{1, 2, \dots, m\}$, оценим разность i -ых компонент векторов Ab и Ab^* :

$$\begin{aligned} |(Ab)_i - (Ab^*)_i| &= |(A(b - b^*))_i| = \left| \sum_{j=1}^n a_{ij}(b_j - b_j^*) \right| \\ &\leq \sqrt{\sum_{j=1}^n a_{ij}^2} \sqrt{\sum_{j=1}^n (b_j - b_j^*)^2} \end{aligned}$$

в силу неравенства Коши-Буняковского. Поэтому $(Ab)_i \rightarrow (Ab^*)_i$ при $b \rightarrow b^*$ для любого номера i .

3.3в Матричные нормы

Помимо векторов основным объектом вычислительной линейной алгебре являются также матрицы. По этой причине нам будут нужны матричные нормы — для того, чтобы оценивать «величину» той или иной матрицы, а также для того, чтобы ввести расстояние между матрицами как

$$\text{dist}(A, B) := \|A - B\|, \quad (3.18)$$

где A, B — вещественные или комплексные матрицы.

Множество матриц само является линейным векторным пространством, а матрица — это составной многомерный объект, в значительной степени аналогичный вектору. Поэтому вполне естественно прежде всего потребовать от матричной нормы тех же свойств, что и для векторной нормы. Формально, матричной нормой на множестве вещественных или комплексных $m \times n$ -матриц называют вещественнозначную функцию $\|\cdot\|$, удовлетворяющую следующим условиям (аксиомам нормы):

(МН1) $\|A\| \geq 0$ для любой матрицы A , причём $\|A\| = 0 \Leftrightarrow A = 0$
 — неотрицательность,

(МН2) $\|\alpha A\| = |\alpha| \cdot \|A\|$ для любых матрицы A и $\alpha \in \mathbb{R}$ или $\alpha \in \mathbb{C}$
 — абсолютная однородность,

(МН3) $\|A + B\| \leq \|A\| + \|B\|$ для любых матриц A, B, C
 — «неравенство треугольника».

Но условия (МН1)–(МН3) выражают взгляд на матрицу, как на «вектор размерности $m \times n$ ». Они явно недостаточным, если мы хотим учесть специфику матриц как объектов, между которыми определена также операция умножения. В частности, множество всех квадратных матриц фиксированного размера наделено более богатой структурой, нежели линейное векторное пространство, и обычно в связи с ним используют уже термин «кольцо» или «алгебра», обозначающее множество с двумя взаимносогласованными бинарными операциями — сложением и умножением (см. [23, 40]). Связь нормы матриц с операцией их умножения отражает четвёртая аксиома матричной нормы:

(МН4) $\|AB\| \leq \|A\| \cdot \|B\|$ для любых матриц A, B
 — «субмультипликативность».⁶

Особую ценность и в теории, и на практике представляют ситуации, когда нормы векторов и нормы матриц, с которыми они совместно рассматриваются и на которые умножаются, существуют не сами по себе, но в некотором смысле согласованы друг с другом. Инструментом такого согласования может как раз-таки выступать аксиома субмультипликативности МН4, понимаемая в расширенном смысле, т. е. для любых матриц A и B таких размеров, что произведение AB имеет смысл. В частности, она должна быть верна для $n \times 1$ -матриц B , являющихся векторами из \mathbb{R}^n .

Определение 3.3.4 *Векторная норма $\|\cdot\|$ и матричная норма $\|\cdot\|'$ называются согласованными, если*

$$\|Ax\| \leq \|A\|' \cdot \|x\| \quad (3.19)$$

для любой матрицы A и всех векторов x .

⁶Приставка «суб-» означает «меньше», «ниже» и т. п. В этом смысле неравенства треугольника ВН3 и МН3 можно называть «субаддитивностью» норм.

Рассмотрим примеры конкретных матричных норм.

Пример 3.3.1 Фробениусова норма матрицы $A = (a_{ij})$ определяется как

$$\|A\|_F = \left(\sum_{i,j} |a_{ij}|^2 \right)^{1/2}.$$

Ясно, что она удовлетворяет первым трём аксиомам матричной нормы просто потому, что задаётся совершенно аналогично евклидовой векторной норме $\|\cdot\|_2$. Для обоснования свойства субмультипликативности рассмотрим

$$\|AB\|_F^2 = \sum_{i,j} \left| \sum_k a_{ik} b_{kj} \right|^2.$$

В силу неравенства Коши-Буняковского (3.16)

$$\left| \sum_k a_{ik} b_{kj} \right|^2 \leq \left(\sum_k a_{ik}^2 \right) \left(\sum_l b_{lj}^2 \right),$$

поэтому

$$\begin{aligned} \|AB\|_F^2 &\leq \sum_{i,j} \left(\sum_k a_{ik}^2 \right) \left(\sum_l b_{lj}^2 \right) \\ &= \sum_{i,j,k,l} a_{ik}^2 b_{lj}^2 = \left(\sum_{i,k} a_{ik}^2 \right) \left(\sum_{l,j} b_{lj}^2 \right) \\ &= \|A\|_F^2 \|B\|_F^2, \end{aligned}$$

что и требовалось.

Если считать, что B — это матрица размера $n \times 1$, т. е. вектор длины n , то выполненные оценки показывают, что фробениусова норма матрицы согласована с евклидовой векторной нормой $\|\cdot\|_2$, с которой она совпадает для векторов. ■

Пример 3.3.2 Матричная норма

$$\|A\|_{\max} = n \max_{i,j} |a_{ij}|,$$

определённая на множестве квадратных $n \times n$ -матриц, является аналогом чебышёвской нормы векторов $\|\cdot\|_\infty$, отличаясь от неё лишь постоянным множителем для матриц фиксированного размера. По этой причине выполнение первых трех аксиом матричной нормы для $\|A\|_{\max}$ очевидно. Но необходимость удовлетворить аксиоме субмультипликативности вызывает появление в выражении для $\|A\|_{\max}$ множителя n перед $\max |a_{ij}|$:

$$\begin{aligned} \|AB\|_{\max} &= n \max_{i,j} \left| \sum_{k=1}^n a_{ik} b_{kj} \right| \leq n \max_{i,j} \left(\sum_{k=1}^n |a_{ik}| |b_{kj}| \right) \\ &\leq n \left(\sum_{k=1}^n \max_{i,k} |a_{ik}| \max_{k,j} |b_{kj}| \right) \\ &\leq n^2 \max_{i,j} |a_{ij}| \max_{i,j} |b_{ij}| = \|A\|_{\max} \|B\|_{\max}. \end{aligned}$$

Ясно, что без этого множителя выписанная выше цепочка неравенств была бы неверной.

Небольшая модификация проведённых выкладок показывает также, что норма $\|A\|_{\max}$ согласована с чебышёвской нормой векторов. Кроме того, несложно устанавливается, что $\|A\|_{\max}$ согласована с евклидовой векторной нормой. ■

В связи с последним примером, следует отметить, что аксиома субмультипликативности МН4 накладывает на матричные нормы более серьёзные ограничения, чем может показаться на первый взгляд. В частности, матричные нормы нельзя произвольно масштабировать, умножая на какое-то число.

Оказывается, среди матричных норм квадратных матриц нет таких, которые не были бы ни с чем согласованными. Иными словами, справедливо

Предложение 3.3.4 *Для любой нормы квадратных матриц можно подобрать подходящую норму векторов, с которой матричная норма будет согласована.*

Доказательство. Для данной нормы $\|\cdot\|'$ на множестве $n \times n$ -матриц определим норму $\|v\|$ для n -вектора v как $\|(v, v, \dots, v)\|'$, т. е. как норму матрицы (v, v, \dots, v) , составленной из n штук векторов v как из

столбцов. Выполнение всех аксиом векторной нормы для $\|v\|$ очевидным образом следует из аналогичных свойств рассматриваемой нормы матрицы.

Опираясь на субмультипликативность матричной нормы, имеем

$$\begin{aligned}\|Av\| &= \|(Av, Av, \dots, Av)\|' = \|A \cdot (v, v, \dots, v)\|' \\ &\leq \|A\|' \cdot \|(v, v, \dots, v)\|' = \|A\|' \cdot \|v\|,\end{aligned}$$

так что требуемое согласование действительно будет достигнуто. ■

3.3г Подчинённые матричные нормы

В предшествующем пункте мы могли видеть, что с заданной векторной нормой могут быть согласованы различные матричные нормы. И наоборот, для матричной нормы возможна согласованность со многими векторными нормами. В этих условиях при проведении различных преобразований и выводе оценок наиболее выгодно оперировать согласованными матричными нормами, которые принимают как можно меньшие значения. Тогда неравенства, получающиеся в результате применения в выкладках соотношения (3.19), будут более точными и позволят получить более тонкие оценки результата. Например, конкретная оценка нормы погрешности может оказать сильное влияние на количество итераций, которые мы вынуждены будем сделать в итерационном численном методе для достижения той или иной точности приближённого решения.

Пусть дана векторная норма $\|\cdot\|$ и зафиксирована матрица A . Из требования согласованности (3.19) вытекает неравенство для согласованной нормы матрицы $\|A\|$:

$$\|A\| \geq \|Ax\|/\|x\|, \quad (3.20)$$

где x — произвольный вектор. Как следствие, значения всех матричных норм от A , согласованных с данной векторной нормой $\|\cdot\|$, ограничены снизу выражением

$$\sup_{x \neq 0} \frac{\|Ax\|}{\|x\|},$$

которое скоро (3.20) должно быть справедливым для любого ненулевого вектора x .

Предложение 3.3.5 Для любой фиксированной векторной нормы $\|\cdot\|$ соотношением

$$\|A\|' = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} \quad (3.21)$$

задаётся матричная норма.

Доказательство. Отметим прежде всего, что в случае конечномерных векторных пространств \mathbb{R}^n и \mathbb{C}^n вместо «sup» в выражении (3.21) можно брать «max». В самом деле,

$$\sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{x \neq 0} \left\| A \frac{x}{\|x\|} \right\| = \sup_{\|y\|=1} \|Ay\|,$$

а задаваемая условием $\|y\| = 1$ единичная сфера любой нормы замкнута и ограничена, т. е. компактна в \mathbb{R}^n или \mathbb{C}^n [40, 50]. Непрерывная функция $\|Ay\|$ достигает на этом компактном множестве своего максимума. Таким образом, в действительности

$$\|A\|' = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|y\|=1} \|Ay\|.$$

Проверим теперь для нашей конструкции выполнение аксиом нормы. Если $A \neq 0$, то найдётся ненулевой вектор y , такой что $Ay \neq 0$. Ясно, что его можно считать нормированным, т. е. $\|y\| = 1$. Тогда $\|Ay\| > 0$, и потому $\max_{\|y\|=1} \|Ay\| > 0$, что доказывает для $\|\cdot\|'$ первую аксиому нормы.

Абсолютная однородность для $\|\cdot\|'$ доказывается тривиально. Покажем для (3.21) справедливость неравенства треугольника. Очевидно,

$$\|(A + B)y\| \leq \|Ay\| + \|By\|,$$

и потому

$$\begin{aligned} \max_{\|y\|=1} \|(A + B)y\| &\leq \max_{\|y\|=1} (\|Ay\| + \|By\|) \\ &\leq \max_{\|y\|=1} \|Ay\| + \max_{\|y\|=1} \|By\|, \end{aligned}$$

что и требовалось.

Приступая к обоснованию субмультипликативности, отметим, что по самому построению $\|Ax\| \leq \|A\|' \|x\|$ для любого вектора x . По этой причине

$$\begin{aligned} \|AB\|' &= \max_{\|y\|=1} \|(AB)y\| = \|ABz\| \quad \text{для некоторого } z \text{ с } \|z\| = 1 \\ &\leq \|A\|' \cdot \|Bz\| \leq \|A\|' \cdot \max_{\|z\|=1} \|Bz\| = \|A\|' \|B\|'. \end{aligned}$$

Это завершает доказательство предложения. ■

Доказанный результат мотивирует

Определение 3.3.5 *Матричная норма, определяемая для заданной векторной нормы $\|\cdot\|$ на линейном векторном пространстве X как*

$$\|A\|' = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|y\|=1} \|Ay\|,$$

называется подчинённой к $\|\cdot\|$ матричной нормой (или индуцированной, или операторной нормой).

Последний термин — операторная норма — популярен потому, что конструкция этой нормы хорошо отражает взгляд на матрицу как на оператор, задающий отображения линейных векторных пространств $\mathbb{R}^n \rightarrow \mathbb{R}^m$ или $\mathbb{C}^n \rightarrow \mathbb{C}^m$. Операторная норма показывает максимальную величину растяжения по норме, которую получает в сравнении с исходным вектором его образ при действии данного оператора.

Несмотря на хорошие свойства подчинённых матричных норм, их определение не отличается большой конструктивностью, так как привлекает операцию взятия максимума. Естественно задаться вопросом о том, существуют ли вообще достаточно простые и обозримые выражения для матричных норм, подчинённых тем или иным векторным нормам. Какими являются подчинённые матричные нормы для рассмотренных выше векторных норм $\|\cdot\|_1$, $\|\cdot\|_2$ и $\|\cdot\|_\infty$? С другой стороны, являются ли матричные нормы $\|A\|_F$ (фробениусова) и $\|A\|_{\max}$ подчинёнными для каких-либо векторных норм?

Ответ на последний вопрос отрицателен. В самом деле, для единичной $n \times n$ -матрицы I имеем

$$\|I\|_{\max} = n, \quad \|I\|_F = \sqrt{n},$$

тогда как из определения подчинённой нормы следует, что должно быть

$$\|I\| = \sup_{\|y\|=1} \|Iy\| = \max_{\|y\|=1} \|y\| = 1. \quad (3.22)$$

Ответом на первые два вопроса является

Предложение 3.3.6 *Для векторной 1-нормы подчинённой матричной нормой для $m \times n$ -матрицы является*

$$\|A\|_1 = \max_{1 \leq j \leq n} \left(\sum_{i=1}^m |a_{ij}| \right)$$

— максимальная сумма модулей элементов по столбцам.

Для чебышёвской векторной нормы (∞ -нормы) подчинённой матричной нормой для $m \times n$ -матрицы является

$$\|A\|_\infty = \max_{1 \leq i \leq m} \left(\sum_{j=1}^n |a_{ij}| \right)$$

— максимальная сумма модулей элементов по строкам.

Матричная норма, подчинённая евклидовой норме векторов $\|x\|_2$, есть $\|A\|_2 = \sigma_{\max}(A)$ — наибольшее сингулярное число матрицы A .

Доказательство. Для обоснования первой части предложения выпишем следующую цепочку преобразований и оценок

$$\begin{aligned} \|Ax\|_1 &= \sum_{i=1}^m |(Ax)_i| = \sum_{i=1}^m \left| \sum_{j=1}^n a_{ij}x_j \right| \leq \sum_{i=1}^m \sum_{j=1}^n |a_{ij}x_j| \\ &= \sum_{i=1}^m \sum_{j=1}^n |a_{ij}| |x_j| = \sum_{j=1}^n \sum_{i=1}^m |a_{ij}| |x_j| = \sum_{j=1}^n \left(|x_j| \sum_{i=1}^m |a_{ij}| \right) \\ &\leq \left(\max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \right) \cdot \sum_{j=1}^n |x_j| = \|A\|_1 \|x\|_1, \end{aligned} \quad (3.23)$$

из которой вытекает $\|A\|_1 \leq \|Ax\|_1 / \|x\|_1$. При этом все неравенства в цепочке (3.23) обращаются в равенства для вектора x в виде столбца единичной $n \times n$ -матрицы с тем номером j , на котором достигается

$\max_j \sum_{i=1}^m |a_{ij}|$. Как следствие, на этом векторе достигается наибольшее значение отношения $\|Ax\|_1/\|x\|_1$ из определения подчинённой матричной нормы.

Аналогичным образом доказывается и вторая часть предложения.

Приступая к обоснованию последней части предложения рассмотрим $n \times n$ -матрицу A^*A . Она является эрмитовой, её собственные числа вещественны и неотрицательны, будучи квадратами сингулярных чисел матрицы A и, возможно, ещё нулями (см. Предложение 3.2.4). Унитарным преобразованием подобия (ортогональным в вещественном случае) матрица A^*A может быть приведена к диагональному виду: $A^*A = U^*AU$, где U — унитарная $n \times n$ -матрица, Λ — диагональная $n \times n$ -матрица, имеющая на диагонали числа σ_i^2 , $i = 1, 2, \dots, \min\{m, n\}$, т. е. квадраты сингулярных чисел σ_i матрицы A , и, возможно, ещё нули в случае $m < n$.

Далее имеем

$$\begin{aligned} \|A\|_2 &= \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \max_{x \neq 0} \frac{\sqrt{x^*A^*Ax}}{\sqrt{x^*x}} = \max_{x \neq 0} \frac{\sqrt{x^*U^*AUx}}{\sqrt{x^*U^*Ux}} \\ &= \max_{x \neq 0} \frac{\sqrt{(Ux)^*\Lambda(Ux)}}{\sqrt{(Ux)^*Ux}} = \max_{y \neq 0} \frac{\sqrt{y^*\Lambda y}}{\sqrt{y^*y}} = \max_{y \neq 0} \sqrt{\frac{\sum_i \sigma_i^2 y_i^2}{\sum_i y_i^2}} \\ &\leq \max_{y \neq 0} \left(\sigma_{\max}(A) \sqrt{\frac{\sum_i y_i^2}{\sum_i y_i^2}} \right) = \sigma_{\max}(A), \end{aligned}$$

где в выкладках применена замена переменных $y = Ux$. Кроме того, полученная для $\|A\|_2$ оценка достижима: достаточно взять в качестве вектора y столбец единичной $n \times n$ -матрицы с номером, равным месту элемента $\sigma_{\max}^2(A)$ на диагонали в Λ , а в самом начале выкладок положить $x = U^*y$. ■

Норму матриц $\|\cdot\|_2$, подчинённую евклидовой векторной норме, часто называют также *спектральной нормой* матриц. Для симметричных матриц она равна наибольшему из модулей собственных чисел и совпадает с так называемым спектральным радиусом матрицы (см. 3.3ж).

Отметим, что спектральная норма матриц не является абсолютной нормой (см. Пример 3.1.3), т. е. она зависит не только от абсолютных значений элементов матрицы. В то же время, $\|\cdot\|_1$ и $\|\cdot\|_\infty$ — это абсолютные матричные нормы, что следует из вида их выражений.

3.3д Топология на множествах матриц

Совершенно аналогично случаю векторов можно рассмотреть топологическую структуру на множестве матриц. Будем говорить, что матричная переменная A сходится к пределу A^* относительно фиксированной нормы матриц (сходится по норме), если $\|A - A^*\| \rightarrow 0$. Матричные нормы назовём *топологически эквивалентными* (или просто *эквивалентными*), если предельный переход в одной норме влечёт существование предела в другой, и обратно. Опять таки, в силу известного факта из математического анализа в конечномерном линейном пространстве всех $m \times n$ -матриц все нормы эквивалентны. Тем не менее, конкретные константы эквивалентности играют огромную роль при выводе различных оценок, и их значения даёт следующее

Предложение 3.3.7 *Для квадратных $n \times n$ -матриц*

$$\begin{aligned} \frac{1}{\sqrt{n}} \|A\|_2 &\leq \|A\|_1 \leq \sqrt{n} \|A\|_2, \\ \frac{1}{\sqrt{n}} \|A\|_\infty &\leq \|A\|_2 \leq \sqrt{n} \|A\|_\infty, \\ \frac{1}{n} \|A\|_1 &\leq \|A\|_\infty \leq n \|A\|_1. \end{aligned}$$

Доказательство. Докажем первое двустороннее неравенство. Имеет место очевидная оценка

$$\|A\|_1 = \max_{\|y\|_1 \leq 1} \|Ay\|_1 \leq \max_{\|y\|_1 \leq 1} (\sqrt{n} \|Ay\|_2)$$

в силу первого неравенства из Предложения 3.3.2. Кроме того, из него следует, что множество векторов y , удовлетворяющих $\|y\|_1 \leq 1$, включается во множество векторов, определяемых условием $\|y\|_2 \leq 1$. По этой причине

$$\max_{\|y\|_1 \leq 1} \|Ay\|_2 \leq \max_{\|y\|_2 \leq 1} \|Ay\|_2 = \|A\|_2,$$

так что в целом действительно $\|A\|_1 \leq \sqrt{n} \|A\|_2$.

С другой стороны, в силу того же первого неравенства из Предложения 3.3.2

$$\|A\|_1 = \max_{\|y\|_1 \leq 1} \|Ay\|_1 \geq \max_{\|y\|_1 \leq 1} \|Ay\|_2.$$

Но множество векторов $\|y\|_1 \leq 1$ не более, чем в \sqrt{n} меньше, чем множество векторов, удовлетворяющих $\|y\|_2 \leq 1$. Как следствие,

$$\max_{\|y\|_1 \leq 1} \|Ay\|_2 \geq \frac{1}{\sqrt{n}} \max_{\|y\|_2 \leq 1} \|Ay\|_2 = \frac{1}{\sqrt{n}} \|A\|_2.$$

Объединяя два выписанных неравенства, получаем требуемое. ■

Как и для векторов, помимо сходимости по норме введём также *поэлементную сходимость* матриц, при которой одна матрица сходится к другой тогда и только тогда, когда все элементы первой матрицы сходятся к соответствующим элементам второй:

$$\begin{aligned} A = (a_{ij}) &\rightarrow A^* = (a_{ij}^*) \text{ в } \mathbb{R}^{m \times n} \text{ или } \mathbb{C}^{m \times n} \\ &\Updownarrow \\ a_{ij} &\rightarrow a_{ij}^* \text{ в } \mathbb{R} \text{ или } \mathbb{C} \text{ для всех индексов } i, j. \end{aligned}$$

Из эквивалентности матричных норм следует, в частности, существование для любой нормы $\|\cdot\|$ такой константы C , что

$$\max_{i,j} |a_{ij}| \leq \max_{1 \leq j \leq n} \left(\sum_{i=1}^m |a_{ij}| \right) = \|A\|_1 \leq C \|A\|$$

(вместо 1-нормы матриц в этой выкладке можно было бы взять, к примеру, ∞ -норму). Поэтому $|a_{ij}| \leq C \|A\|$, так что сходимость последовательности матриц в любой норме равносильна сходимости последовательностей всех элементов этих матриц.

В целом множество матриц с введённым на нём посредством (3.18) расстоянием для любой матричной нормы является полным метрическим пространством, т. е. любая фундаментальная («сходящаяся в себе») последовательность имеет в нём предел. Это следует из предшествующего рассуждения и из факта полноты вещественной оси \mathbb{R} и комплексной плоскости \mathbb{C} .

В заключение этой темы отметим, что в вычислительной линейной алгебре понятия норм векторов и матриц впервые стали широко использоваться В.Н. Фаддеевой в монографии [81], которая предшествовала капитальной книге [44] и вошла в неё составной частью.

3.3е Энергетическая норма

Ещё одной важной и популярной конструкцией нормы является так называемая энергетическая норма векторов, которая порождается какой-либо симметричной положительно-определённой матрицей (эрмитовой в комплексном случае). Если A — такая матрица, то выражение $\langle Ax, y \rangle$, как нетрудно проверить, есть симметричная билинейная положительно-определённая форма, т. е. скалярное произведение векторов x и y . Следовательно, можно определить норму вектора x , как

$$\|x\|_A := \sqrt{\langle Ax, x \rangle},$$

т. е. как корень из произведения x на себя в этом новом скалярном произведении. Она называется *энергетической нормой* вектора x относительно матрицы A , и нижний индекс указывает на эту порождающую матрицу. Её часто называют также A -нормой векторов, если в задаче имеется в виду какая-то конкретная симметричная положительно-определённая матрица A . Термин «энергетическая» происходит из-за аналогии выражения для этой нормы с выражениями для различных видов энергии (см. также §3.10а).

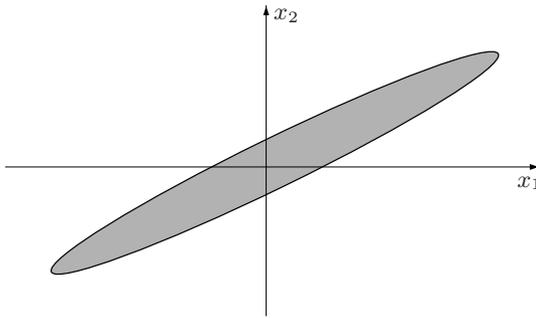


Рис. 3.7. Шар единичного радиуса в энергетической норме при значительном разбросе спектра порождающей матрицы

Так как симметричная матрица может быть приведена к диагональному виду ортогональными преобразованиями подобия, то

$$A = Q^T D Q,$$

где Q — ортогональная матрица, $D = \text{diag} \{ \lambda_1, \lambda_2, \dots, \lambda_n \}$ — диагональная матрица, на главной диагонали которой стоят собственные

значения λ_i матрицы A . Поэтому

$$\begin{aligned} \|x\|_A &= \sqrt{\langle Ax, x \rangle} = \sqrt{\langle Q^T D Q x, x \rangle} \\ &= \sqrt{\langle D Q x, Q x \rangle} = \sqrt{\langle D y, y \rangle} = \left(\sum_i \lambda_i^2 y_i^2 \right)^{1/2}. \end{aligned} \quad (3.24)$$

где $y = Qx$. Таким образом, в системе координат, которая получается из исходной ортогональным преобразованием $x = Q^T y$, линии уровня энергетической нормы, т. е. поверхности $\|x\|_A = \text{const}$, являются эллипсоидами. Они тем более вытянуты, чем больше различаются между собой λ_i , т. е. чем больше разброс собственных чисел матрицы A .

Из сказанного вытекает характерная особенность энергетической нормы, которая в ряде случаев оборачивается её недостатком, — возможность существенного искажения обычного геометрического масштаба объектов по разным направлениям (своеобразная анизотропия). Она вызывается разбросом собственных значений порождающей матрицы A и приводит к тому, что векторы из \mathbb{R}^n , имеющие одинаковую энергетическую норму, существенно различны по обычной евклидовой длине, и наоборот (Рис. 3.7). С другой стороны, использование энергетической нормы, которая порождена матрицей, фигурирующей в постановке задачи (системе линейных алгебраических уравнений, задаче на собственные значения и т. п.) часто является удобным и оправданным, а альтернативы ему очень ограничены. Примеры будут рассмотрены в §3.10б, §3.10в и §3.10г.

Из общего факта эквивалентности любых норм в конечномерном линейном пространстве следует, что энергетическая норма эквивалентна рассмотренным выше матричным нормам $\|\cdot\|_1$, $\|\cdot\|_2$, $\|\cdot\|_\infty$, фробениусовой норме и норме $\|\cdot\|_{\max}$. Но интересно знать конкретные константы эквивалентности. Из выражения (3.24) следует, что

$$\left(\min_i |\lambda_i| \right) \|x\|_2 \leq \|x\|_A \leq \left(\max_i |\lambda_i| \right) \|x\|_2.$$

Другие двусторонние неравенства для энергетической нормы можно получить на основе Предложения 3.3.7.

Выражения для матричных норм, которые подчинены энергетической норме векторов или просто согласованы с нею, выписываются сложно и даже не всегда могут быть указаны в явном и несложно вычисляемом виде. Тем не менее, можно привести полезный и красивый

результат на эту тему, который будет далее использован при исследовании метода наискорейшего спуска в §3.10б:

Предложение 3.3.8 *Если S — матрица, которая является значением некоторого многочлена от матрицы A , порождающей энергетическую норму $\|\cdot\|_A$, то для любого вектора x справедливо*

$$\|Sx\|_A \leq \|S\|_2 \|x\|_A. \quad (3.25)$$

Доказательство. Воспользуемся спектральным разложением матрицы A , представив её в виде $A = QDQ^*$, где Q — ортогональная (унитарная в общем случае) матрица, а D — диагональная матрица с положительными собственными значениями A по диагонали. Ясно, что матрица S — симметричная (эрмитова) одновременно с A , причём для неё справедливо аналогичное разложение $S = Q\Sigma Q^*$ с той же самой матрицей Q , где $\Sigma = \text{diag}\{s_1, s_2, \dots, s_n\}$ — диагональная матрица, имеющая по диагонали собственные числа S . Тогда $S^* = Q\Sigma Q^*$ и потому

$$\begin{aligned} \|Sx\|_A^2 &= \langle ASx, Sx \rangle = \langle A Q \Sigma Q^* x, Q \Sigma Q^* x \rangle \\ &= \langle Q D Q^* Q \Sigma Q^* x, Q \Sigma Q^* x \rangle = \langle D \Sigma Q^* x, \Sigma Q^* x \rangle \\ &= \langle \Sigma D \Sigma Q^* x, Q^* x \rangle = \langle \Sigma^2 D Q^* x, Q^* x \rangle \\ &\leq s_1^2 \langle D Q^* x, Q^* x \rangle = s_1^2 \langle Q^* D Q^* x, x \rangle \\ &= \|S\|_2^2 \langle Ax, x \rangle = \|S\|_2^2 \|x\|_A^2, \end{aligned}$$

где s_1 — наибольшее сингулярное число матрицы S , т. е. $s_1 = \|S\|_2$. ■

3.3ж Спектральный радиус

Определение 3.3.6 *Спектральным радиусом квадратной матрицы называется наибольший из модулей её собственных чисел.*

Эквивалентное определение: спектральным радиусом матрицы называется наименьший из радиусов кругов комплексной плоскости \mathbb{C} с центром в начале координат, который содержит весь спектр матрицы. Эта трактовка хорошо объясняет и сам термин. Обычно спектральный радиус матрицы A обозначают $\rho(A)$.

Спектральный радиус матрицы — неотрицательное число, которое в общем случае может не совпадать ни с одним из собственных значений (см. Рис. 3.8). Но если матрица неотрицательна, т. е. все её элементы — неотрицательные вещественные числа, то наибольшее по модулю собственное значение такой матрицы также неотрицательно и, таким образом, равно спектральному радиусу матрицы. Кроме того, неотрицательным может быть выбран соответствующий собственный вектор. Эти утверждения составляют содержание теоремы Перрона-Фробениуса, одного из главных результатов теории неотрицательных матриц (см. [9, 35, 50]).

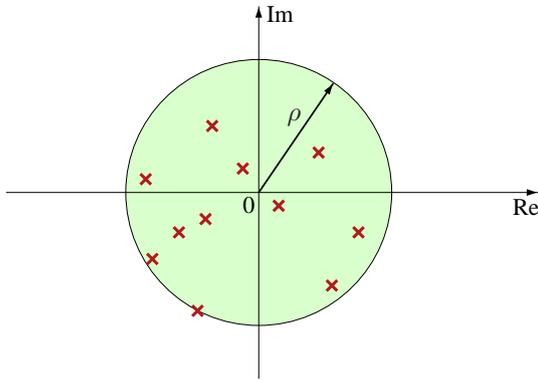


Рис. 3.8. Иллюстрация спектрального радиуса матрицы: крестиками обозначены точки спектра.

Предложение 3.3.9 *Спектральный радиус матрицы не превосходит любой её нормы.*

Доказательство. Рассмотрим сначала случай, когда матрица является комплексной.

Пусть λ — собственное значение матрицы A , а $v \neq 0$ — соответствующий собственный вектор, так что $Av = \lambda v$. Воспользуемся тем установленным в §3.3в фактом (Предложение 3.3.4), что любая матричная норма согласована с некоторой векторной нормой, и возьмём от обеих частей равенства $Av = \lambda v$ норму, согласованную с рассматриваемой

$\|A\|$. Получим

$$\|A\| \cdot \|v\| \geq \|Av\| = \|\lambda v\| = |\lambda| \cdot \|v\|, \quad (3.26)$$

где $\|v\| > 0$, и потому сокращение на эту величину обеих частей неравенства (3.26) даёт $\|A\| \geq |\lambda|$. Коль скоро наше рассуждение справедливо для любого собственного значения λ , то в самом деле $\max \lambda = \rho(A) \leq \|A\|$.

Рассмотрим теперь случай вещественной $n \times n$ -матрицы A . Если λ — её вещественное собственное значение, то проведённые выше рассуждения остаются полностью справедливыми. Если же λ — комплексное собственное значение матрицы A , то комплексным является и соответствующий собственный вектор v . Тогда цепочку соотношений (3.26) выписать нельзя, поскольку согласованная векторная норма определена лишь для вещественных векторов из \mathbb{R}^n .

Выполним *комплексификацию* рассматриваемого линейного пространства, т. е. вложим его в более широкое линейное векторное пространство над полем комплексных чисел. В формальных терминах мы переходим от \mathbb{R}^n к пространству $\mathbb{R}^n \oplus i\mathbb{R}^n$, где i — мнимая единица (т. е. скаляр, обладающий свойством $i^2 = -1$), $i\mathbb{R}^n$ — это множество всех произведений iy для $y \in \mathbb{R}^n$, а « \oplus » означает прямую сумму линейных пространств (см. [10, 23, 35, 83]).

Элементами $\mathbb{R}^n \oplus i\mathbb{R}^n$ служат упорядоченные пары $(x, y)^\top$, где $x, y \in \mathbb{R}^n$. Сложение и умножение на скаляр $(\alpha + i\beta) \in \mathbb{C}$ определяются для них следующим образом

$$(x, y)^\top + (x', y')^\top = (x + x', y + y')^\top, \quad (3.27)$$

$$(\alpha + i\beta) \cdot (x, y)^\top = (\alpha x - \beta y, \alpha y + \beta x)^\top. \quad (3.28)$$

Введённые пары векторов $(x, y)^\top$ обычно записывают в виде $x + iy$, причём x и y называются соответственно вещественной и мнимой частями вектора из $\mathbb{R}^n \oplus i\mathbb{R}^n$. Линейный оператор, действующий на $\mathbb{R}^n \oplus i\mathbb{R}^n$ и продолжающий линейное отображение на \mathbb{R}^n , порождаемое матрицей A , может быть представлен в матричном виде как

$$\mathcal{A} = \begin{pmatrix} A & 0 \\ 0 & A \end{pmatrix}. \quad (3.29)$$

Его блочно-диагональный вид объясняется тем, что согласно формуле (3.28) для любого $\alpha \in \mathbb{R}$

$$\alpha \cdot (x, y)^\top = (\alpha x, \alpha y)^\top,$$

и потому вещественная матрица A независимо действует на вещественную и мнимую части векторов из построенного комплексного пространства $\mathbb{R}^n \oplus i\mathbb{R}^n$.

Без какого-либо ограничения общности можно считать, что рассматриваемая нами норма матрицы, т.е. $\|A\|$, является подчинённой (операторной) нормой, так как такие нормы являются наименьшими из всех согласованных матричных норм (см. §3.3г). Если предложение будет обосновано для подчинённых матричных норм, то оно тем более будет верным для всех прочих норм матриц.

Пусть $\|\cdot\|$ — векторная норма в \mathbb{R}^n , которой подчинена наша матричная норма. Зададим в $\mathbb{R}^n \oplus i\mathbb{R}^n$ норму векторов как $\|(x, y)^T\| = \|x\| + \|y\|$. Тогда ввиду (3.29) и с помощью рассуждений, аналогичных доказательству Предложения 3.3.6, нетрудно показать, что подчинённая матричная норма для \mathcal{A} во множестве $2n \times 2n$ -матриц есть $\|\mathcal{A}\| = \max\{\|A\|, \|A\|\} = \|A\|$. Кроме того, теперь для \mathcal{A} справедливы рассуждения о связи нормы и спектрального радиуса, проведённые в начале доказательства для случая комплексной матрицы, т.е.

$$\rho(A) = \rho(\mathcal{A}) \leq \|\mathcal{A}\| = \|A\|.$$

Это и требовалось доказать. ■

Для симметричных и эрмитовых матриц спектральный радиус есть норма, которая совпадает со спектральной матричной нормой $\|\cdot\|_2$. Это следует из Предложения 3.3.6 и того факта, что для симметричных и эрмитовых матриц сингулярные числа равны абсолютным значениям собственных чисел. Но для матриц общего вида спектральный радиус матричной нормой не является. Хотя для любого скаляра α справедливо

$$\rho(\alpha A) = |\alpha| \rho(A),$$

т.е. спектральный радиус обладает абсолютной однородностью, аксиома неотрицательности матричной нормы (МН1) и неравенство треугольника (МН3) для него не выполняются.

Во-первых, для ненулевой матрицы

$$\begin{pmatrix} 0 & 1 & & 0 \\ & 0 & 1 & \\ & & \ddots & \ddots \\ \mathbf{0} & & & 0 & 1 \\ & & & & & 0 \end{pmatrix} \quad (3.30)$$

— жордановой клетки, отвечающей собственному значению 0, спектральный радиус равен нулю. Во-вторых, если A — матрица вида (3.30), то $\rho(A^\top) = \rho(A) = 0$, но $\rho(A + A^\top) > 0$. Это вытекает из того, что симметричная матрица $A + A^\top$ — ненулевая, поэтому $\|A + A^\top\|_2 > 0$ и, как следствие, наибольший из модулей её собственных значений строго больше нуля. Получается, что неверно «неравенство треугольника»

$$\rho(A + A^\top) \leq \rho(A) + \rho(A^\top).$$

Тем не менее, спектральный радиус является важной характеристикой матрицы, которая описывает асимптотическое поведение её степеней.

Предложение 3.3.10 Пусть A — квадратная матрица, вещественная или комплексная. Для сходимости степеней $A^k \rightarrow 0$ при $k \rightarrow \infty$ необходимо, чтобы $\rho(A) < 1$, т. е. чтобы спектральный радиус матрицы A был меньше 1.

Доказательство. Пусть λ — собственное число матрицы A (возможно, комплексное), а $v \neq 0$ — соответствующий ему собственный вектор (который также может быть комплексным). Тогда $Av = \lambda v$, и потому

$$\begin{aligned} A^2v &= A(Av) = A(\lambda v) = \lambda(Av) = \lambda^2v, \\ A^3v &= A(A^2v) = A(\lambda^2v) = \lambda^2(Av) = \lambda^3v, \\ \dots & \quad \dots \quad , \end{aligned}$$

так что в целом

$$(A^k)v = (\lambda^k)v. \tag{3.31}$$

Если последовательность степеней A^k , $k = 0, 1, 2, \dots$, сходится к нулевой матрице, то при фиксированном векторе v нулевой предел имеет и вся левая часть выписанного равенства. Как следствие, к нулевому вектору должна сходиться и правая часть в (3.31), причём $v \neq 0$. Это возможно лишь в случае $|\lambda| < 1$. ■

Ниже в §3.9б мы увидим, что условие $\rho(A) < 1$ является, в действительности, также достаточным для сходимости к нулю степеней матрицы A .

Рассуждения, с помощью которых доказано Предложение 3.3.10, можно продолжить и несложно вывести весьма тонкие свойства спек-

трального радиуса. Возьмём от обеих частей равенства (3.31) какую-нибудь векторную норму:

$$\|A^k v\| = \|\lambda^k v\|.$$

Поэтому $\|A^k\| \|v\| \geq |\lambda^k| \|v\|$ для согласованной матричной нормы $\|A\|$, так что после сокращения на $\|v\| \neq 0$ получаем

$$\|A^k\| \geq |\lambda|^k \quad \text{для всех } k = 0, 1, 2, \dots$$

По этой причине для любого собственного значения матрицы имеет место оценка

$$|\lambda| \leq \inf_{k \in \mathbb{N}} \|A^k\|^{1/k},$$

или, иными словами,

$$\rho(A) \leq \inf_{k \in \mathbb{N}} \|A^k\|^{1/k}. \quad (3.32)$$

Так как всякая матричная норма всегда согласована с какой-то векторной, то выведенное неравенство справедливо для любой матричной нормы. Оно является обобщением Предложения 3.3.9, переходя в него при $k = 1$.

Уточнением неравенства (3.32) является *формула Гельфанда*

$$\rho(A) = \lim_{k \rightarrow \infty} \|A^k\|^{1/k},$$

которая также верна для любой из матричных норм. Её доказательство можно найти, к примеру, в [50].

Предложение 3.3.10, неравенство (3.32) и формула Гельфанда, показывают, что с помощью спектрального радиуса адекватно описывается асимптотическое поведение норм степеней матрицы. Несмотря на то, что матрица является сложным составным объектом, нормы её степеней ведут себя примерно так же, как геометрическая прогрессия со знаменателем, равным спектральному радиусу. Например, для матрицы (3.30) или любой ей подобной n -ая степень зануляется, и это свойство обнаруживается спектральным радиусом.

3.3з Матричный ряд Неймана

Как известно из математического анализа, операцию суммирования можно обобщить на случай бесконечного числа слагаемых, и такие бес-

конечные суммы называются *рядами*. При этом *суммой ряда* называется предел (если он существует) сумм конечного числа слагаемых, когда количество слагаемых неограниченно возрастает. Совершенно аналогичная конструкция применима также к суммированию векторов и матриц, а не только чисел. Именно, суммой матричного ряда

$$\sum_{k=0}^{\infty} A^{(k)},$$

где $A^{(k)}$, $k = 0, 1, 2, \dots$, — матрицы одного размера, мы будем называть предел частичных сумм $\sum_{k=0}^N A^{(k)}$ при $N \rightarrow \infty$. В этом определении $A^{(k)}$ могут быть и векторами.

Предложение 3.3.11 Пусть X — квадратная матрица и $\|X\| < 1$ в некоторой матричной норме. Тогда матрица $(I - X)$ неособенна, для обратной матрицы справедливо представление

$$(I - X)^{-1} = \sum_{k=0}^{\infty} X^k, \quad (3.33)$$

и имеет место оценка

$$\|(I - X)^{-1}\| \leq \frac{1}{1 - \|X\|}. \quad (3.34)$$

Фигурирующий в правой части равенства (3.33) аналог геометрической прогрессии для матриц называется *матричным рядом Неймана*.

Доказательство. Покажем неособенность матрицы $(I - X)$. Если это не так, то $(I - X)v = 0$ для некоторого ненулевого вектора v . Тогда $Xv = v$, и, беря от обеих частей этого равенства векторную норму, согласованную с матричной нормой, в которой $\|X\| < 1$ по условию Предложения, мы получим

$$\|X\| \|v\| \geq \|Xv\| = \|v\|.$$

В случае, когда $v \neq 0$, можем сократить обе части полученного неравенства на положительную величину $\|v\|$, что даёт $\|X\| \geq 1$. Следовательно, при условии $\|X\| < 1$ и ненулевых v равенство $(I - X)v = 0$ невозможно.

Обозначим $S_N = \sum_{k=0}^N X^k$ — частичную сумму матричного ряда Неймана. Коль скоро

$$\begin{aligned} \|S_{N+p} - S_N\| &= \left\| \sum_{k=N+1}^{N+p} X^k \right\| \leq \sum_{k=N+1}^{N+p} \|X^k\| \leq \sum_{k=N+1}^{N+p} \|X\|^k \\ &= \|X\|^{N+1} \cdot \frac{1 - \|X\|^p}{1 - \|X\|} \rightarrow 0 \end{aligned}$$

при $N \rightarrow \infty$ и любых целых положительных p , то последовательность S_N является фундаментальной (последовательностью Коши) в полном метрическом пространстве квадратных матриц с расстоянием, порождённым рассматриваемой нормой $\|\cdot\|$. Следовательно, частичные суммы S_N ряда Неймана имеют предел $S = \lim_{N \rightarrow \infty} S_N$, причём

$$(I - X)S_N = (I - X)(I + X + X^2 + \dots + X^N) = I - X^{N+1} \rightarrow I$$

при $N \rightarrow \infty$, поскольку тогда $\|X^{N+1}\| \leq \|X\|^{N+1} \rightarrow 0$. Так как этот предел S удовлетворяет соотношению $(I - X)S = I$, можем заключить, что $S = (I - X)^{-1}$.

Наконец,

$$\|(I - X)^{-1}\| = \left\| \sum_{k=0}^{\infty} X^k \right\| \leq \sum_{k=0}^{\infty} \|X^k\| \leq \sum_{k=0}^{\infty} \|X\|^k = \frac{1}{1 - \|X\|},$$

где для бесконечных сумм неравенство треугольника может быть обосновано предельным переходом по аналогичным неравенствам для конечных сумм. Это завершает доказательство Предложения. ■

Матричный ряд Неймана является простейшим из матричных степенных рядов, т. е. сумм вида

$$\sum_{k=0}^{\infty} c_k X^k$$

где X — квадратная матрица и c_k , $k = 0, 1, 2, \dots$, — счётный набор коэффициентов. С помощью матричных степенных рядов можно определять значения аналитических функций от матриц (например, экспоненту, логарифм, синус, косинус и т. п. от матрицы), просто подставляя

матрицу вместо аргумента в степенные разложения для соответствующих функций. Эта важная и интересная тема, находящая многочисленные приложения, развивается в рамках так называемой теории представлений линейных операторов.

3.4 Приложения сингулярного разложения

3.4а Исследование особенностей и ранга матриц

Рассмотренное в §3.2д сингулярное разложение матрицы может служить основой для вычислительных технологий решения некоторых важных математических задач. Рассмотрим первой задачу об определении того, особенна или неособенна матрица.

Квадратная диагональная матрица неособенна тогда и только тогда, когда все её диагональные элементы не равны нулю. Из сингулярного разложения матрицы следует, что произвольная квадратная матрица неособенна тогда и только тогда, когда её сингулярные числа — ненулевые.

Хотя определение особенности или неособенности матрицы обычно ассоциируется с исследованием определителя этой матрицы, наиболее надёжным в вычислительном отношении способом проверки особенности/неособенности является исследование сингулярных чисел матрицы. Хотя эта процедура более трудоёмка, чем нахождение определителя, она гораздо более предпочтительна в силу существенно большей устойчивости к ошибкам. Кроме того, величина ненулевого определителя матрицы является неадекватным признаком того, насколько близка матрица к особенной: с помощью умножения матрицы на подходящее число значение её определителя можно сделать любым, тогда как мера линейной независимости столбцов матрицы или её строк при этом никак не изменится.

Рассмотрим теперь задачу о вычислении ранга матрицы. Согласно определению, ранг — это количество линейно независимых вектор-строк или вектор-столбцов матрицы, с помощью которых можно линейным комбинированием породить всю матрицу. Фактически, ранг — это число независимых параметров, задающих матрицу. В таком виде хорошо видна важность ранга в задачах обработки данных, когда нам необходимо выявить какие-то закономерности в массивах данных, полученных в результате наблюдений или опытов. С помощью ранга

можно увидеть, к примеру, что все данные суть линейные комбинации немногих порождающих.

Ранг матрицы не зависит непрерывно от её элементов. Выражаясь языком, который развивается в Главе 4 (§4.2), можно сказать, что задача вычисления ранга матрицы не является вычислительно-корректной. Как следствие, совершенно точное определение ранга в условиях «зашумлённых» данных, которые искажены случайными помехами и ошибками измерений, не имеет смысла. Нам нужно, как правило, знать «приближённый ранг», и при прочих равных условиях для его нахождения более предпочтителен тот метод, который менее чувствителен к ошибкам и возмущениям в данных. Под «приближённым рангом» естественно понимать ранг матрицы, приближённо равной исходной в смысле некоторой разумной нормы.

Ясно, что ранг диагональной матрицы равен числу её ненулевых диагональных элементов. Ортогональные преобразования сохраняют линейную независимость. Таким образом, ранг любой матрицы равен количеству её ненулевых сингулярных чисел, что следует из сингулярного разложения (3.12), т. е. представления

$$A = U\Sigma V^*$$

Другой способ нахождения ранга матрицы может состоять в приведении её к так называемому строчно-ступенчатому виду с помощью преобразований, которые использовались в прямом ходе метода Гаусса. Но в условиях неточных данных и неточных арифметических операций на ЭВМ строчно-ступенчатая форма является очень ненадёжным инструментом. Использование сингулярного разложения — более надёжный и достаточно эффективный подход к нахождению ранга матрицы.

3.46 Решение систем линейных уравнений

Если для матрицы A известно сингулярное разложение (3.12), то система линейных алгебраических уравнений $Ax = b$ может быть переписана эквивалентным образом как

$$U\Sigma V^*x = b.$$

Отсюда решение легко находится в виде

$$x = V\Sigma^{-1}U^*b.$$

Получается, что для вычисления решения мы должны умножить вектор правой части на ортогональную матрицу, затем разделить компоненты результата на сингулярные числа и, наконец, ещё раз умножить получившийся вектор на другую ортогональную матрицу. Вычислительной работы здесь существенно больше, чем при реализации, к примеру, метода исключения Гаусса (см. §3.6б) или других прямых методов решения СЛАУ, особенно с учётом того, что сингулярное разложение матрицы системы нужно ещё найти. Но описанный путь безупречен с вычислительной точки зрения, так как позволяет без накопления ошибок найти решение системы и, кроме того, проанализировать состояние её разрешимости.

Напомним, что с геометрической точки зрения преобразования, осуществляемые ортогональными матрицами, являются обобщениями поворотов и отражений: они сохраняют длины и углы. Поэтому в вычислительном отношении умножения на ортогональные матрицы обладают очень хорошими свойствами, так как не увеличивают ошибок округлений и других погрешностей. Ниже в §3.5б мы взглянём на этот факт с другой стороны. Отличие в поведении и результатах метода Гаусса и метода, основанного на сингулярном разложении, особенно зримо в случае, когда матрица системы «почти особенна».

3.4в Малоранговые приближения матрицы

Пусть A — $m \times n$ -матрица, u_k и v_k — это её k -ые нормированные левый и правый сингулярные векторы, а Υ_k обозначает их внешнее произведение, т. е.

$$\Upsilon_k = u_k v_k^*. \quad (3.35)$$

Отметим, что Υ_k — $m \times n$ -матрица ранга 1. Тогда сингулярное разложение (3.12) матрицы A равносильно её представлению в виде суммы

$$A = \sum_{k=1}^n \sigma_k \Upsilon_k,$$

где σ_i , $i = 1, 2, \dots, \min\{m, n\}$, — сингулярные числа матрицы A . Если $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$, и мы «обрубаем» выписанную сумму после p -го слагаемого, то получающаяся матрица называется p -ранговым приближением данной матрицы:

$$A_p = \sum_{k=1}^p \sigma_k \Upsilon_k, \quad (3.36)$$

Это в самом деле матрица ранга p , что следует из её сингулярного представления, а погрешность, с которой она приближает исходную матрицу, равна

$$\sum_{k=p+1}^n \sigma_k \Upsilon_k.$$

Величина этой погрешности решающим образом зависит от величины сингулярных чисел $\sigma_{p+1}, \dots, \sigma_{\min\{m,n\}}$, соответствующих отброшенным слагаемым в (3.35). Более точно, погрешность p -рангового приближения характеризуется следующим замечательным свойством:

Теорема 3.4.1 Пусть σ_k, u_k и v_k — сингулярные числа и левые и правые сингулярные векторы $m \times n$ -матрицы A соответственно. Если $p < n$ и

$$A_p = \sum_{k=1}^p \sigma_k u_k v_k^*$$

— p -ранговое приближение матрицы A , то

$$\|A - A_p\|_2 = \min_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq p}} \|A - B\|_2 = \sigma_{p+1}.$$

Иными словами, относительно спектральной нормы p -ранговое приближение матрицы обеспечивает наименьшее отклонение от первоначальной матрицы среди всех матриц ранга не более p .

Доказательство. Предположим, что найдётся такая матрица B , имеющая ранг $\text{rank}(B) \leq p$, что $\|A - B\|_2 < \|A - A_p\|_2 = \sigma_{p+1}$. Тогда существует $(n-p)$ -мерное подпространство $W \subset \mathbb{C}^n$, для которого справедливо $w \in W \Rightarrow Bw = 0$. При этом для любого $w \in W$ мы имеем $Aw = (A - B)w$, так что

$$\|Aw\|_2 = \|(A - B)w\|_2 \leq \|A - B\|_2 \|w\|_2 < \sigma_{p+1} \|w\|_2.$$

Таким образом, W является $(n-p)$ -мерным подпространством в \mathbb{C}^n , в котором $\|Aw\|_2 < \sigma_{p+1} \|w\|_2$.

Но в \mathbb{C}^n имеется $(p+1)$ -мерное подпространство, образованное векторами v , для которых $\|Av\|_2 \geq \sigma_{p+1} \|v\|_2$. Это подпространство, являющееся линейной оболочкой первых $p+1$ правых сингулярных векторов матрицы A . Поскольку сумма размерностей этого подпространства и

подпространства W превосходит n , размерности всего пространства, то должен существовать ненулевой вектор, лежащий в них обоих. Это приводит к противоречию. ■

Совершенно аналогичный результат справедлив для фробениусовой нормы матриц, и исторически он был обнаружен даже раньше, чем Теорема 3.4.1:

Теорема 3.4.2 (теорема Экарта-Янга [87]) Пусть σ_k , u_k и v_k — сингулярные числа и левые и правые сингулярные векторы $m \times n$ -матрицы A соответственно. Если $p < n$ и

$$A_p = \sum_{k=1}^p \sigma_k u_k v_k^*$$

— p -ранговое приближение матрицы A , то

$$\|A - A_p\|_F = \min_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq p}} \|A - B\|_F = \sigma_{p+1},$$

где $\|\cdot\|_F$ — фробениусова норма матриц. Иными словами, относительно фробениусовой нормы p -ранговое приближение матрицы обеспечивает наименьшее отклонение от первоначальной матрицы среди всех матриц ранга не более p .

Доказательство опускается.

Итак, если младшие сингулярные числа матрицы достаточно малы, то вместо неё можно взять p -ранговое приближение вида (3.36). Оно более «экономно», с меньшим числом параметров, представляет исходную матрицу.

3.4г Метод главных компонент

В качестве важного и интересного практического примера, который иллюстрирует понятия ранга матрицы, матричной нормы, сингулярных чисел и сингулярных векторов матрицы и пр. рассмотрим так называемый *метод главных компонент*, широко применяемый в анализе данных и статистике.

Во многих практических задачах приходится иметь дело с большими массивами числовых данных, характеризующих какой-либо объект

или явление. Предположим для определённости, что набор параметров (свойств, признаков и т. п.) рассматриваемого объекта характеризуется вектор-строкой из n чисел, и мы имеем m штук таких векторов, относящихся, к примеру, к отдельным измерениям. Полученные данные образуют вещественную $m \times n$ -матрицу, которую мы обозначим через A .

Нередко возникает необходимость сжатия данных, т. е. уменьшения числа n параметров объекта с тем, чтобы оставшиеся p признаков, $p \leq n$, всё-таки «достаточно наиболее полно» описывали всю совокупность накопленной об объекте информации, содержащейся в матрице A . Метод главных компонент является одним из способов решения поставленного вопроса, который в формализованном виде принимает следующую форму: существует ли в \mathbb{R}^n ортонормированный базис $\{e_1, e_2, \dots, e_p\}$, $p < n$, в котором рассматриваемые нами данные, содержащиеся в матрице A , будут представлены в наиболее экономичной (удобной, красивой и т. п.) форме?

В качестве меры «близости» матриц мы можем брать различные расстояния, получая различные постановки задач. Одним из практически наиболее важных является расстояние, порождённое фробениусовой нормой матриц, которое имеет ясный вероятностно-статистический смысл: оно совпадает с так называемой выборочной дисперсией набора данных. Для фробениусовой нормы матриц математическая задача ставится следующим образом. Нужно найти такой ортонормированный базис $\{e_1, e_2, \dots, e_p\}$ в \mathbb{R}^n , $p \leq n$, что квадратичное отклонение исходных векторов данных $A_i = (a_{i1}, a_{i2}, \dots, a_{in})^\top$ от их приближений $X^{(i)} = \sum_{j=1}^p x_{ij}e_j$ в этом базисе было бы наименьшим возможным для всех $i = 1, 2, \dots, m$.

Описанная выше процедура обработки матрицы данных называется *методом главных компонент* в случае применения к матрице данных, которая центрирована путём вычитания из каждого столбца его среднего значения. При этом *компонентами* называются правые сингулярные векторы v_k , а масштабированные левые сингулярные векторы $\sigma_k u_k$ носят название *долей*. Метод главных компонент обычно описывают в терминах собственных чисел и собственных векторов так называемой ковариационной матрицы $A^\top A$, но подход, основанный на сингулярном разложении, часто лучше с вычислительной точки зрения.

Другая ситуация, в которой часто прибегают к методу главных компонент и которая не связана с необходимостью сжатия данных, — это желание выделить из данных наиболее значимые *факторы*, т. е. комби-

нации переменных, наиболее существенные для рассматриваемого объекта или явления. Здесь и пригождается понятие ранга матрицы или же приближённого ранга для случая неточных данных.

Приведённая выше теорема Экарта-Янга даёт математическую основу для решения поставленной задачи. Следует отметить, что соответствующие результаты неоднократно переоткрывались статистиками и, по-видимому, первым метод главных компонент предложил К. Пирсон в начале XX века, который отметил, что искомый минимум достигается в том случае, если базис $\{e_1, e_2, \dots, e_p\}$ берётся в виде собственных векторов так называемой ковариационной матрицы $C = A^T A$, отвечающих её p наибольшим собственным значениям. На современном языке можно сказать, что искомый базис составлен из старших сингулярных векторов матрицы данных A , а «главными компонентами» обычно именуют компоненты разложения векторов данных по этому базису.

3.5 Обусловленность систем линейных уравнений

3.5а Число обусловленности матриц

В этом параграфе общие идеи и понятия, развитые в §1.3, рассматриваются в приложении к задаче решения системы линейных уравнений. В частности, мы вводим количественную меру чувствительности решения по отношению к вариациям матрицы и вектора правой части.

Рассмотрим систему линейных алгебраических уравнений

$$Ax = b$$

с неособенной квадратной матрицей A и вектором правой части $b \neq 0$, а также систему

$$(A + \Delta A) \tilde{x} = b + \Delta b,$$

где $\Delta A \in \mathbb{R}^{n \times n}$ и $\Delta b \in \mathbb{R}^n$ — возмущения матрицы и вектора правой части. Насколько сильно ненулевое решение \tilde{x} возмущённой системы может отличаться от решения x исходной системы уравнений?

Пусть это отличие есть $\Delta x = \tilde{x} - x$, так что $\tilde{x} = x + \Delta x$, и потому

$$(A + \Delta A)(x + \Delta x) = b + \Delta b.$$

Вычитая из этого равенства исходную невозмущённую систему уравнений, получим

$$(\Delta A)x + (A + \Delta A)\Delta x = \Delta b, \quad (3.37)$$

или

$$(\Delta A)(x + \Delta x) + A\Delta x = \Delta b,$$

так что

$$\Delta x = A^{-1}(-(\Delta A)\tilde{x} + \Delta b).$$

Для оценки величины изменения решения Δx воспользуемся какой-нибудь удобной векторной нормой. Применяя её к обеим частям полученного соотношения, будем иметь

$$\|\Delta x\| \leq \|A^{-1}\| \cdot (\|\Delta A\| \|\tilde{x}\| + \|\Delta b\|)$$

при согласовании используемых векторных и матричных норм. Предполагая, что возмущённое решение \tilde{x} не равно нулю, можем поделить обе части на $\|\tilde{x}\| > 0$, придя к неравенству

$$\begin{aligned} \frac{\|\Delta x\|}{\|\tilde{x}\|} &\leq \|A^{-1}\| \cdot \left(\|\Delta A\| + \frac{\|\Delta b\|}{\|\tilde{x}\|} \right) \\ &= \|A^{-1}\| \|A\| \cdot \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|A\| \cdot \|\tilde{x}\|} \right). \end{aligned} \quad (3.38)$$

Это весьма практичная *апостериорная оценка* относительной погрешности решения, которую удобно применять после того, как приближённое решение системы уже найдено.⁷ Коль скоро $\|A\| \cdot \|\tilde{x}\| \geq \|A\tilde{x}\| \approx \|b\|$, то знаменатель второго слагаемого в скобках из правой части неравенства «приблизительно не превосходит» $\|b\|$. Поэтому полученной оценке (3.38) путём некоторого огрубления можно придать более элегантный вид

$$\frac{\|\Delta x\|}{\|\tilde{x}\|} \approx \leq \|A^{-1}\| \|A\| \cdot \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|b\|} \right), \quad (3.39)$$

в котором справа задействованы относительные погрешности в матрице A и правой части b .

⁷От латинского словосочетания «a posteriori», означающего знание, полученное из опыта. Под «опытом» здесь понимается процесс решения задачи.

Фигурирующая в оценках (3.38) и (3.39) величина $\|A^{-1}\| \|A\|$, на которую суммарно умножаются ошибки в матрице и правой части, имеет своё собственное название, так как играет важнейшую роль в вычислительной линейной алгебре.

Определение 3.5.1 Для квадратной неособенной матрицы A величина $\|A^{-1}\| \|A\|$ называется её числом обусловленности (относительно выбранной нормы матриц).

Мы будем обозначать число обусловленности матрицы A посредством $\text{cond}(A)$, иногда с индексом, указывающим выбор нормы.⁸ Если же матрица A особенна, то удобно положить $\text{cond}(A) = +\infty$. Это соглашение оправдывается тем, что обычно $\|A^{-1}\|$ неограниченно возрастает при приближении матрицы A к множеству особенных матриц.

Выведем теперь *априорную* оценку относительной погрешности ненулевого решения, которая не будет опираться на знание вычисленного решения и годится для получения оценки *до* решения СЛАУ.⁹

После вычитания точного уравнения из приближённого мы получили (3.37):

$$(\Delta A)x + (A + \Delta A)\Delta x = \Delta b.$$

Отсюда

$$\begin{aligned} \Delta x &= (A + \Delta A)^{-1}(-(\Delta A)x + \Delta b) \\ &= (A(I + A^{-1}\Delta A))^{-1}(-(\Delta A)x + \Delta b) \\ &= (I + A^{-1}\Delta A)^{-1}A^{-1}(-(\Delta A)x + \Delta b). \end{aligned}$$

Беря интересующую нас векторную норму от обеих частей этого равенства и пользуясь далее условием согласования с матричной нормой, субмультипликативностью и неравенством треугольника, получим

$$\|\Delta x\| \leq \|(I + A^{-1}\Delta A)^{-1}\| \cdot \|A^{-1}\| \cdot (\|\Delta A\| \|x\| + \|\Delta b\|),$$

откуда после деления обеих частей на $\|x\| > 0$:

$$\frac{\|\Delta x\|}{\|x\|} \leq \|(I + A^{-1}\Delta A)^{-1}\| \cdot \|A^{-1}\| \cdot \left(\|\Delta A\| + \frac{\|\Delta b\|}{\|x\|} \right).$$

⁸В математической литературе для числа обусловленности матрицы A можно встретить также обозначения $\mu(A)$ или $\kappa(A)$.

⁹От латинского словосочетания «a priori», означающего в философии знание, полученное до опыта и независимо от него.

Предположим, что возмущение ΔA матрицы A не слишком велико, так что выполнено условие

$$\|\Delta A\| \leq \frac{1}{\|A\|}.$$

Тогда

$$\|A^{-1}\Delta A\| \leq \|A^{-1}\| \|\Delta A\| < 1,$$

и обратная матрица $(I + A^{-1}\Delta A)^{-1}$ разлагается в матричный ряд Неймана (3.33). Соответственно, мы можем воспользоваться вытекающей из этого оценкой (3.34). Тогда

$$\begin{aligned} \frac{\|\Delta x\|}{\|x\|} &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\Delta A\|} \cdot \left(\|\Delta A\| + \frac{\|\Delta b\|}{\|x\|} \right) \\ &= \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \|\Delta A\|} \cdot \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|A\| \|x\|} \right) \\ &\leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \cdot \frac{\|\Delta A\|}{\|A\|}} \cdot \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|b\|} \right), \end{aligned} \quad (3.40)$$

поскольку $\|A\| \|x\| \geq \|Ax\| = \|b\|$.

Оценка (3.40) — важная априорная оценка относительной погрешности численного решения системы линейных алгебраических уравнений через оценки относительных погрешностей её матрицы и правой части. Если величина $\|\Delta A\|$ достаточно мала, то множитель усиления относительной ошибки в данных

$$\frac{\text{cond}(A)}{1 - \text{cond}(A) \cdot \frac{\|\Delta A\|}{\|A\|}}$$

близок к числу обусловленности матрицы A .

Понятие числа обусловленности матрицы и полученные с его помощью оценки имеют большое теоретическое значение, но их практическая полезность напрямую зависит от наличия эффективных способов вычисления или хотя бы приближённого оценивания числа обусловленности матриц. Фактически, определение числа обусловленности требует знания некоторых характеристик обратной матрицы, и в случае общих матричных норм хорошего решения задачи оценивания $\text{cond}(A)$ не существует до сих пор.

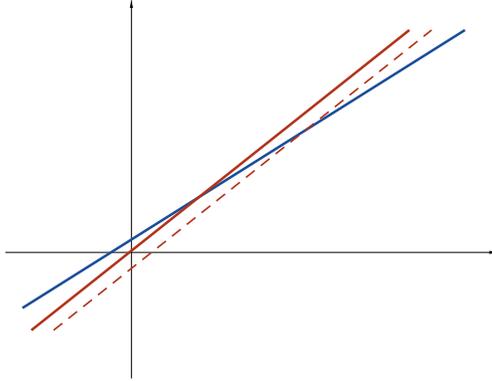


Рис. 3.9. Иллюстрация возмущения решения системы линейных уравнений с плохой обусловленностью матрицы.

Тем не менее, существует практически важный частный случай, когда нахождение числа обусловленности матрицы может быть вполне достаточно эффективно. Это случай спектральной матричной нормы $\|\cdot\|_2$, подчинённой евклидовой норме векторов.

Напомним (Предложение 3.2.5), что для любой неособенной квадратной матрицы A справедливо равенство $\sigma_{\max}(A^{-1}) = \sigma_{\min}^{-1}(A)$, и поэтому относительно спектральной нормы число обусловленности матрицы есть

$$\text{cond}_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}.$$

Этот результат помогает понять большую роль сингулярных чисел в современной вычислительной линейной алгебре и важность алгоритмов для их нахождения. В совокупности с ясным геометрическим смыслом евклидовой векторной нормы (2-нормы) это вызывает преимущественное использование этих норм для многих задач теории и практики.

Наконец, если матрица A симметрична (эрмитова), то её сингулярные числа совпадают с модулями собственных значений, и тогда

$$\text{cond}_2(A) = \frac{\max_i |\lambda_i(A)|}{\min_i |\lambda_i(A)|} \quad (3.41)$$

— спектральное число обусловленности равно отношению наибольшего и наименьшего по модулю собственных значений матрицы. Для сим-

метричных положительно определённых матриц эта формула принимает совсем простой вид

$$\text{cond}_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}.$$

3.56 Примеры хорошообусловленных и плохообусловленных матриц

Условимся называть матрицу *хорошо обусловленной*, если её число обусловленности невелико. Напротив, если число обусловленности матрицы велико, станем говорить, что матрица *плохо обусловлена*. Естественно, что эти определения имеют неформальный характер, так как зависят от нестрогих понятий «невелико» и «велико». Тем не менее, они весьма полезны в практическом отношении, в частности, потому, что позволяют сделать наш язык более выразительным.

Отметим, что для любой подчинённой матричной нормы

$$\text{cond}(A) = \|A^{-1}\| \|A\| \geq \|A^{-1}A\| = \|I\| = 1$$

в силу (3.22), и поэтому соответствующее число обусловленности матрицы всегда не меньше единицы. Для произвольных матричных норм полученное неравенство тем более верно в силу того, что подчинённые нормы принимают наименьшие значения.

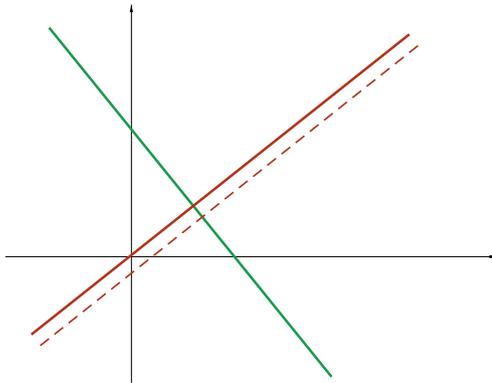


Рис. 3.10. Иллюстрация возмущения решения системы линейных уравнений с хорошей обусловленностью матрицы.

Примером матриц, обладающих наилучшей возможной обусловленностью относительно спектральной нормы, являются ортогональные матрицы, для которых $\text{cond}_2(Q) = 1$. Действительно, если Q ортогональна, то $\|Qx\|_2 = \|x\|_2$ для любого вектора x . Следовательно, $\|Q\|_2 = 1$. Кроме того, $Q^{-1} = Q^T$ и тоже ортогональна, а потому $\|Q^{-1}\|_2 = 1$.

Самым популярным содержательным примером плохообусловленных матриц являются, пожалуй, матрицы Гильберта $H_n = (h_{ij})$, которые встретились нам в §2.10г при обсуждении среднеквадратичного приближения алгебраическими полиномами на интервале $[0, 1]$. Это симметричные матрицы, образованные элементами

$$h_{ij} = \frac{1}{i+j-1}, \quad i, j = 1, 2, \dots, n,$$

так что, к примеру,

$$H_3 = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{pmatrix}.$$

Число обусловленности матриц Гильберта исключительно быстро растёт в зависимости от их размера n . Воспользовавшись какими-либо стандартными процедурами для вычисления числа обусловленности матриц (встроенными, к примеру, в системы компьютерной математики Scilab, МАТЛАВ, Octave, Maple и им подобные), нетрудно найти следующие числовые данные:

$$\text{cond}_2(H_2) = 19.3,$$

$$\text{cond}_2(H_3) = 524,$$

...

$$\text{cond}_2(H_{10}) = 1.6 \cdot 10^{13},$$

...

Существует общая формула [97, 98] :

$$\text{cond}_2(H_n) = O\left(\frac{(1+\sqrt{2})^{4n}}{\sqrt{n}}\right) \approx O(34^n/\sqrt{n}),$$

где O — «о большое», известный из математического анализа символ Э. Ландау (см. стр. 95). Интересно, что матрицы, обратные к матрицам

Гильберта могут быть вычислены аналитически в явном виде [88]. Они имеют целочисленные элементы, которые также очень быстро растут с размерностью.

На этом фоне для матрицы Вандермонда (2.7) оценка снизу для числа обусловленности (см. [54])

$$\text{cond}_2(V(x_0, x_1, \dots, x_n)) \geq \sqrt{2} \frac{(1 + \sqrt{2})^{n-1}}{\sqrt{n+1}}$$

представляется существенно более скромной.¹⁰ Но она и не хороша, так что матрицы Вандермонда можно называть «умеренно плохообусловленными».

3.5в Практическое применение числа обусловленности матриц

Оценки (3.38) и (3.40) на возмущения решений систем линейных алгебраических уравнений являются неулучшаемыми на всём множестве матриц, векторов правых частей и их возмущений. Более точно, для данной матрицы эти оценки достигаются на каких-то векторах правой части и возмущениях матрицы и правой части. Но «плохая обусловленность» матрицы не всегда означает высокую чувствительность решения конкретной системы по отношению к тем или иным конкретным возмущениям. Если, к примеру, правая часть имеет нулевые компоненты в направлении сингулярных векторов, отвечающих наименьшим сингулярным числам матрицы системы, то решение СЛАУ зависит от возмущений этой правой части гораздо слабее, чем показывает оценка (3.40) для спектральной нормы (см. рассуждения в §3.4). И определение того, какова конкретно правая часть по отношению к матрице СЛАУ — плохая или не очень — не менее трудно, чем само решение данной системы линейных уравнений.

Из сказанного должна вытекать известная осторожность и осмотрительность по отношению к выводам, которые делаются о практической разрешимости и достоверности решений какой-либо системы линейных уравнений лишь на основании того, велико или мало число обусловленности их матрицы. Тривиальный пример: число обусловленности диагональной матрицы может быть сколь угодно большим, но решение СЛАУ с такими матрицами почти никаких проблем не вызывает!

¹⁰ Аналогичные по смыслу, но более слабые экспоненциальные оценки снизу для числа обусловленности матрицы Вандермонда выводятся также в книге [41].

Наконец, число обусловленности малоприспособно для оценки разброса решения СЛАУ при значительных и больших изменениях элементов матрицы и правой части (начиная с нескольких процентов от исходного значения). Получаемые при этом с помощью оценок (3.38) и (3.40) результаты типично завышены во много раз (иногда на порядки), и для решения упомянутой задачи более предпочтительны методы интервального анализа (см., к примеру, [84, 93]).

Пример 3.5.1 Рассмотрим 2×2 -систему линейных уравнений

$$\begin{pmatrix} 3 & -1 \\ 0 & 3 \end{pmatrix} x = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

в которой элементы матрицы и правой части заданы неточно, с абсолютной погрешностью 1, так что в действительности можно было бы записать эту систему в неформальном виде как

$$\begin{pmatrix} 3 \pm 1 & -1 \pm 1 \\ 0 \pm 1 & 3 \pm 1 \end{pmatrix} x = \begin{pmatrix} 0 \pm 1 \\ 1 \pm 1 \end{pmatrix}.$$

Фактически, мы имеем совокупность эквивалентных по точности систем линейных уравнений

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} x = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix},$$

у которых элементы матрицы и правой части могут принимать значения из интервалов

$$\begin{array}{lll} a_{11} \in [2, 4], & a_{12} \in [-2, 0], & b_1 \in [-1, 1] \\ a_{12} \in [-1, 1], & a_{22} \in [2, 4], & b_2 \in [0, 2]. \end{array}$$

При этом обычно говорят [84, 93], что задана *интервальная система линейных алгебраических уравнений*

$$\begin{pmatrix} [2, 4] & [-2, 0] \\ [-1, 1] & [2, 4] \end{pmatrix} x = \begin{pmatrix} [-1, 1] \\ [0, 2] \end{pmatrix}. \quad (3.42)$$

Её *множеством решений* называют множество, образованное всевозможными решениями систем линейных алгебраических уравнений той

же структуры, коэффициенты матрицы и компоненты правой части которой принадлежат заданным интервалам. В частности, множество решений рассматриваемой нами системы (3.42) изображено на Рис. 3.11. Мы более подробно рассматриваем интервальные линейные системы уравнений в §4.6.

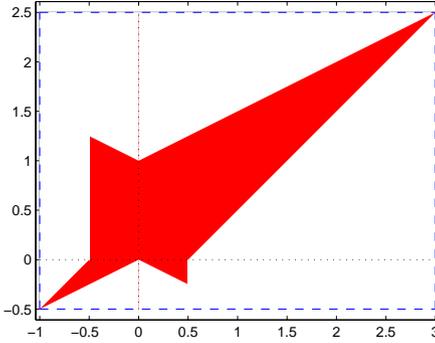


Рис. 3.11. Множество решений интервальной линейной системы (3.42).

Подсчитаем оценки возмущений, которые получаются на основе числа обусловленности для решения системы (3.42). Её можно рассматривать, как систему, получающуюся путём возмущения «средней системы»

$$\begin{pmatrix} 3 & -1 \\ 0 & 3 \end{pmatrix} x = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

на величину

$$\Delta A = \begin{pmatrix} \Delta a_{11} & \Delta a_{12} \\ \Delta a_{21} & \Delta a_{22} \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, \quad \|\Delta A\|_{\infty} \leq 2,$$

в матрице и величину

$$\Delta b = \begin{pmatrix} \Delta b_1 \\ \Delta b_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \|\Delta b\|_{\infty} \leq 1,$$

в правой части. Чебышёвская векторная норма (∞ -норма) используется здесь для оценки Δb потому, что она наиболее адекватно (без искажения формы) описывает возмущение правой части b . Соответствующая ∞ -норма для матрицы ΔA , подчинённая векторной ∞ -норме,

также наиболее уместна в этой ситуации, поскольку обеспечивает наиболее аккуратное согласование вычисляемых оценок.

Обусловленность средней матрицы относительно ∞ -нормы равна 1.778, ∞ -норма средней матрицы равна 4, а ∞ -норма средней правой части — это 1. Следовательно, по формуле (3.40) получаем

$$\frac{\|\Delta x\|}{\|x\|} \lesssim 24.$$

Поскольку решение средней системы есть $\tilde{x} = (\frac{1}{3}, \frac{1}{9})^\top$, и оно имеет ∞ -норму $\frac{1}{3}$, то оценкой разброса решений рассматриваемой системы уравнений является $\tilde{x} \pm \Delta x$, где $\|\Delta x\|_\infty \leq 8$, т. е. двумерный брус¹¹

$$\begin{pmatrix} [-7.667, 8.333] \\ [-7.889, 8.111] \end{pmatrix}.$$

По размерам он в более чем в 4 (четыре) раза превосходит оптимальные (точные) покоординатные оценки множества решений, равные

$$\begin{pmatrix} [-1, 3] \\ [-0.5, 2.5] \end{pmatrix}.$$

Этот брус выделен пунктиром на Рис. 3.11.

При использовании других норм результаты, даваемые формулой (3.40), совершенно аналогичны своей грубостью оценивания возмущений решений. ■

Отметим в заключение этой темы, что задача оценивания разброса решений СЛАУ при вариациях входных данных является в общем случае NP-трудной [90, 91]. Иными словами, если мы не накладываем ограничений на величину возмущений в данных, она требует для своего решения экспоненциально больших трудозатрат.

¹¹Читатель может проверить числовые данные этого примера в любой системе компьютерной математики: Scilab, МАТЛАВ, Octave и т. п.

3.6 Прямые методы решения систем линейных алгебраических уравнений

Решение систем линейных алгебраических уравнений вида

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1, \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2, \\ \vdots \quad \quad \quad \ddots \quad \quad \quad \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n, \end{cases} \quad (3.43)$$

с коэффициентами a_{ij} и свободными членами b_i , или, в краткой форме,

$$Ax = b \quad (3.44)$$

с $m \times n$ -матрицей $A = (a_{ij})$ и m -вектором правой части $b = (b_i)$, является важной математической задачей. Она часто встречается как сама по себе, так и в качестве составного элемента в технологической цепочке решения более сложных задач. Например, решение нелинейных уравнений или систем уравнений часто сводится к последовательности решений линейных уравнений (метод Ньютона).

Следует отметить, что системы линейных алгебраических уравнений не всегда предъявляются к решению в каноническом виде (3.43). Это придаёт дополнительную специфику процессу решения подобных задач и иногда диктует выбор тех или иных методов решения.

Пример 3.6.1 Пусть задана двумерная область $\mathcal{D} = [\underline{x}_1, \bar{x}_1] \times [\underline{x}_2, \bar{x}_2]$, имеющая форму прямоугольника со сторонами, параллельными координатным осям. Рассмотрим в ней численное решение дифференциального уравнения Лапласа

$$\frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2} = 0. \quad (3.45)$$

Уравнением Лапласа описывается, к примеру, распределение температуры стационарного теплового поля, потенциал электростатического поля или же потенциальное течение несжимаемой жидкости. Для определения конкретного решения этого уравнения задают ещё какие-либо краевые условия на границе расчётной области. Мы будем считать за-

данными значения искомой функции $u(x_1, x_2)$ на границе прямоугольника:

$$u(\underline{x}_1, x_2) = \underline{f}(x_2), \quad u(\bar{x}_1, x_2) = \bar{f}(x_2), \quad (3.46)$$

$$u(x_1, \underline{x}_2) = \underline{g}(x_1), \quad u(x_1, \bar{x}_2) = \bar{g}(x_1). \quad (3.47)$$

Рассматриваемую задачу определения функции $u(x_1, x_2)$, которая удовлетворяет уравнению (3.45) внутри области и условиям (3.46)–(3.47) на границе, называют *задачей Дирихле* для уравнения Лапласа.

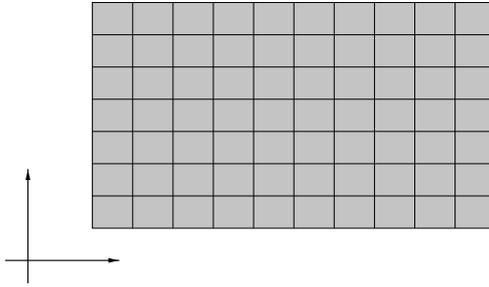


Рис. 3.12. Расчётная область и сетка для численного решения уравнения Лапласа.

Станем решать задачу (3.45)–(3.47) с помощью *конечно-разностного метода*, в котором искомая функция заменяется своим дискретным аналогом, а производные в решаемом уравнении заменяются на разностные отношения. Введём на области \mathcal{D} равномерную прямоугольную сетку, и вместо функции $u(x_1, x_2)$ непрерывного аргумента будем рассматривать её значения в узлах построенной сетки (см. 3.6).

Если обозначить через u_{ij} значение искомой функции u в точке x_{ij} , то после замены вторых производных формулами (2.65) получим следующую систему соотношений

$$\frac{u_{i-1,j} - 2u_{ij} + u_{i+1,j}}{h_1^2} + \frac{u_{i,j-1} - 2u_{ij} + u_{i,j+1}}{h_2^2} = 0, \quad (3.48)$$

$i = 1, 2, \dots, m - 1, j = 1, 2, \dots, n - 1$, для внутренних узлов расчётной

области. На границе области имеем условия

$$u_{i0} = \underline{f}_i, \quad u_{in} = \overline{f}_i, \quad (3.49)$$

$$u_{0j} = \underline{g}_j, \quad u_{mj} = \overline{g}_j, \quad (3.50)$$

$i = 1, 2, \dots, m-1, j = 1, 2, \dots, n-1$. Соотношения (3.48) и (3.46)–(3.47) образуют, очевидно, систему линейных алгебраических уравнений относительно неизвестных u_{ij} , $i = 1, 2, \dots, m-1, j = 1, 2, \dots, n-1$, но она не имеет канонический вид (3.43), так как неизвестные имеют двойные индексы. Конкретный вид (3.43), который получит решаемая система уравнений, зависит от способа выбора базиса в пространстве векторов неизвестных, в частности, от способа перенумерации этих неизвестных, при котором мы образуем из них вектор с одним индексом.

Ясно, что рассмотренный пример может быть сделан ещё более выразительным в трёхмерном случае, когда нам необходимо численно решать трёхмерное уравнение Лапласа. ■

Системы линейных алгебраических уравнений, аналогичные рассмотренной в Примере 3.6.1, где матрица и вектор неизвестных не заданы в явном виде, будем называть системами в *операторной форме*. Не все из изложенных ниже методов решения СЛАУ могут быть непосредственно применены к системам подобного вида.

По характеру вычислительного алгоритма методы решения уравнений и систем уравнений традиционно разделяют на *прямые* и *итерационные*. В прямых методах искомое решение получается в результате выполнения конечной последовательности действий, так что эти методы нередко называют ещё *конечными* или даже *точными*. Напротив, в итерационных методах решение достигается как предел некоторой последовательности приближений, которая конструируется по решаемой системе уравнений.

Одна из основных идей, лежащих в основе прямых методов для решения систем линейных алгебраических уравнений, состоит в том, чтобы эквивалентными преобразованиями привести решаемую систему к наиболее простому виду, из которого решение находится уже непосредственно. В качестве простейших могут выступать системы с диагональными, двухдиагональными, треугольными и т. п. матрицами. Чем меньше ненулевых элементов остаётся в матрице преобразованной системы, тем проще и устойчивее её решение, но, с другой стороны, тем сложнее и неустойчивее приведение к такому виду. На практике обыч-

но стремятся к компромиссу между этими взаимно противоположными требованиями, и в зависимости от целей, преследуемых при решении СЛАУ, приводят её к диагональному (метод Гаусса-Йордана), двухдиагональному (см., к примеру, [65]) или треугольному виду. Мы, в основном, рассмотрим методы, основанные на приведении к треугольному виду.

Наконец, для простоты мы далее подробно разбираем случай систем уравнений (3.43)–(3.44), в которых число неизвестных n равно числу уравнений m , т. е. имеющих квадратную $n \times n$ -матрицу коэффициентов.

3.6а Решение треугольных линейных систем

Напомним, что *треугольными матрицами* называют матрицы, у которых все элементы ниже главной диагонали либо все элементы выше главной диагонали нулевые (так что и нулевые, и ненулевые элементы образуют треугольники):

$$U = \begin{pmatrix} \times & \times & \cdots & \times & \times \\ & \times & \ddots & \times & \times \\ & & \ddots & \vdots & \vdots \\ & \mathbf{0} & & \times & \times \\ & & & & \times \end{pmatrix}, \quad L = \begin{pmatrix} \times & & & & \mathbf{0} \\ \times & \times & & & \\ \times & \times & \ddots & & \\ \vdots & \vdots & \ddots & \times & \\ \times & \times & \cdots & \times & \times \end{pmatrix},$$

где крестиками « \times » обозначены ненулевые элементы. В первом случае говорят о *верхней* (или *правой*) треугольной матрице, а во втором — о *нижней* (или *левой*) треугольной матрице. Соответственно, *треугольными* называются системы линейных алгебраических уравнений, матрицы которых имеют треугольный вид — верхний или нижний.

Рассмотрим для определённости линейную систему уравнений

$$Lx = b \tag{3.51}$$

с неособенной нижней треугольной матрицей $L = (l_{ij})$, так что $l_{ij} = 0$ при $j > i$ и $l_{ii} \neq 0$ для всех $i = 1, 2, \dots, n$. Её первое уравнение содержит только одну неизвестную переменную x_1 , второе уравнение содержит две неизвестных переменных x_1 и x_2 , и т. д., так что в i -е уравнение входят лишь переменные x_1, x_2, \dots, x_i . Найдём из первого уравнения значение x_1 и подставим его во второе уравнение системы, в котором в результате останется всего одна неизвестная переменная x_2 . Вычислим

x_2 и затем подставим известные значения x_1 и x_2 в третье уравнение, из которого определится x_3 . И так далее.

Решение линейной системы (3.51) с нижней треугольной $n \times n$ -матрицей выполняется по следующему простому алгоритму

DO FOR $i = 1$ TO n

$$x_i \leftarrow \left(b_i - \sum_{j < i} l_{ij} x_j \right) / l_{ii}, \quad (3.52)$$

END DO

позволяющему последовательно друг за другом вычислить искомые значения неизвестных переменных, начиная с первой. Этот процесс называется *прямой подстановкой*, коль скоро он выполняется по возрастанию индексов компонент вектора x и его главным содержанием является подстановка, на очередном шаге, уже найденных значений неизвестных в следующее уравнение.

Для решения верхних треугольных систем линейных уравнений существует аналогичный процесс, который называется *обратной подстановкой* — он идёт в обратном направлении, т. е. от x_n к x_1 . Мы рассмотрим его подробнее в следующем разделе.

3.6б Метод Гаусса для решения линейных систем уравнений

Описываемый в этом разделе *метод Гаусса* для решения систем линейных алгебраических уравнений впервые в новом времени был описан К.Ф. Гауссом в 1849 году, хотя письменные источники свидетельствуют о том, что он был известен как минимум за 250 лет до нашей эры.

Хорошо известно, что умножение какого-либо уравнения системы на ненулевое число, а также замена уравнения на его сумму с другим уравнением системы приводят к равносильной системе уравнений, т. е. имеющей те же самые решения. Воспользуемся этими свойствами для преобразования решаемой системы линейных алгебраических уравнений к более простому виду.

Пусть дана система линейных алгебраических уравнений

$$\left\{ \begin{array}{l} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1, \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2, \\ \vdots \quad \quad \quad \ddots \quad \quad \quad \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_m, \end{array} \right.$$

и относительно её коэффициента a_{11} будем предполагать, что $a_{11} \neq 0$. Умножим первое уравнение системы на $(-a_{21}/a_{11})$ и сложим со вторым уравнением. В результате коэффициент a_{21} во втором уравнении занулится, а получившаяся система будет совершенно равносильна исходной.

Проделаем подобное преобразование с остальными — 3-м, 4-м и т. д. до n -го уравнениями системы, т. е. будем умножать первое уравнение на $(-a_{i1}/a_{11})$ и складывать с i -ым уравнением системы. В результате получим равносильную исходной систему линейных алгебраических уравнений, в которой неизвестная переменная x_1 присутствует лишь в первом уравнении. Матрица получившейся СЛАУ станет выглядеть следующим образом:

$$\left(\begin{array}{c|cccc} \times & \times & \times & \cdots & \times \\ \hline 0 & \times & \times & \cdots & \times \\ 0 & \times & \times & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \times & \times & \cdots & \times \end{array} \right),$$

где посредством « \times » обозначены элементы, возможно, не равные нулю.

Рассмотрим теперь в преобразованной системе уравнения со 2-го по n -е. Они образуют квадратную $(n-1) \times (n-1)$ -систему линейных уравнений, в которой неизвестная переменная x_1 уже не присутствует и которую можно решать отдельно, никак не обращаясь к первому уравнению исходной системы. Если элемент на месте (2, 2) не сделался равным нулю, к этой системе можно заново применить вышеописанную процедуру исключения неизвестных. Её результатом будет обнуление поддиагональных элементов 2-го столбца матрицы СЛАУ. И так далее.

Выполнив $(n-1)$ шагов подобного процесса — для 1-го, 2-го, ..., $(n-1)$ -го столбцов матрицы данной системы, мы получим, в конце

концов, линейную систему с верхней треугольной матрицей, которая несложно решается с помощью обратной подстановки, рассмотренной выше в §3.6а. Описанное преобразование системы линейных алгебраических уравнений к равносильному треугольному виду называется *прямым ходом* метода Гаусса, и его псевдокод выглядит следующим образом:

```

DO FOR  $j = 1$  TO  $n - 1$ 
  DO FOR  $i = j + 1$  TO  $n$ 
     $r_{ij} \leftarrow (-a_{ij}/a_{jj})$ 
    DO FOR  $k = j + 1$  TO  $n$ 
       $a_{ik} \leftarrow a_{ik} + r_{ij}a_{jk}$ 
    END DO
     $b_i \leftarrow b_i + r_{ij}b_j$ 
  END DO
END DO

```

(3.53)

Он выражает процесс последовательного обнуления поддиагональных элементов j -го столбца матрицы системы, $j = 1, 2, \dots, n - 1$, и соответствующие преобразования вектора правой части. Матрица системы при этом приводится к верхнему треугольному виду. Далее следует *обратный ход* метода Гаусса для решения полученной верхней треугольной системы, и он является процессом обратной подстановки из §3.6а:

```

DO FOR  $i = n$  DOWNT0 1
  
$$x_i \leftarrow \left( b_i - \sum_{j>i} a_{ij}x_j \right) / a_{ii}$$

END DO

```

(3.54)

Он позволяет последовательно вычислить, в обратном порядке, искомые значения неизвестных, начиная с n -ой. Отметим, что в псевдокоде

(3.53) прямого хода метода Гаусса зануление поддиагональных элементов первых столбцов уже учтено нижней границей внутреннего цикла по k , которая равна $j + 1$.

Помимо изложенной выше вычислительной схемы существует много других версий метода Гаусса. Весьма популярной является, к примеру, *схема единственного деления*. При выполнении её прямого хода сначала делят первое уравнение системы на $a_{11} \neq 0$, что даёт

$$x_1 + \frac{a_{12}}{a_{11}} x_2 + \cdots + \frac{a_{1n}}{a_{11}} x_n = \frac{b_1}{a_{11}}. \quad (3.55)$$

Умножая затем уравнение (3.55) на a_{i1} и вычитая результат из i -го уравнения системы для $i = 2, 3, \dots, n$, добиваются обнуления поддиагональных элементов первого столбца. Затем процедура повторяется в отношении 2-го уравнения и 2-го столбца получившейся СЛАУ, и так далее. Обратный ход совпадает с (3.54).

Схема единственного деления совершенно эквивалентна алгоритму (3.53) и отличается от него лишь тем, что для каждого столбца деление в ней выполняется действительно только один раз, тогда как все остальные операции — это умножение и сложение. С другой стороны, уравнения преобразуемой системы в схеме единственного деления дополнительно масштабируются диагональными коэффициентами при неизвестных, и в некоторых случаях это бывает нежелательно.

3.6в Матричная интерпретация метода Гаусса

Умножение первого уравнения системы на $r_{i1} = -a_{i1}/a_{11}$ и сложение его с i -ым уравнением могут быть представлены в матричном виде как умножение обеих частей системы $Ax = b$ слева на матрицу

$$\begin{pmatrix} 1 & & & & & \\ 0 & 1 & & & & \\ \vdots & & \ddots & & & \\ r_{i1} & & & 1 & & \\ \vdots & & & & 1 & \\ 0 & & 0 & & & 1 \end{pmatrix},$$

которая отличается от единичной матрицы наличием одного дополнительного ненулевого элемента r_{i1} на месте $(i, 1)$. Исключение поддиа-

гональных элементов первого столбца матрицы СЛАУ — это последовательное домножение обеих частей этой системы слева на матрицы

$$\begin{pmatrix} 1 & & & \mathbf{0} \\ r_{21} & 1 & & \\ 0 & & \ddots & \\ \vdots & \mathbf{0} & & 1 \\ 0 & & & 1 \end{pmatrix}, \quad \begin{pmatrix} 1 & & & \mathbf{0} \\ 0 & 1 & & \\ r_{31} & 0 & \ddots & \\ \vdots & \mathbf{0} & & 1 \\ 0 & & & 1 \end{pmatrix},$$

и так далее до

$$\begin{pmatrix} 1 & & & \mathbf{0} \\ 0 & 1 & & \\ \vdots & & \ddots & \\ 0 & \mathbf{0} & & 1 \\ r_{n1} & & & 1 \end{pmatrix}.$$

Нетрудно убедиться, что умножение матриц выписанного выше вида выполняется по простому правилу:

$$\begin{pmatrix} 1 & & & \mathbf{0} \\ & 1 & & \\ r_{i1} & & \ddots & \\ & & & 1 \\ \mathbf{0} & & & \ddots \\ & & & & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & & & \mathbf{0} \\ & 1 & & \\ & & \ddots & \\ r_{k1} & & & 1 \\ \mathbf{0} & & & \ddots \\ & & & & 1 \end{pmatrix} = \begin{pmatrix} 1 & & & \mathbf{0} \\ & 1 & & \\ r_{i1} & & \ddots & \\ & & & 1 \\ r_{k1} & & & \ddots \\ \mathbf{0} & & & & 1 \end{pmatrix}.$$

Оно также остаётся верным в случае, когда у матриц-сомножителей на взаимнодополнительных местах в первом столбце присутствует более одного ненулевого элемента. Следовательно, обнуление поддиагональных элементов первого столбца и соответствующие преобразования правой части в методе Гаусса — это не что иное, как умножение обеих частей СЛАУ слева на матрицу

$$E_1 = \begin{pmatrix} 1 & & & \mathbf{0} & & \\ r_{21} & 1 & & & & \\ r_{31} & 0 & 1 & & & \\ \vdots & & & \ddots & & \\ r_{n1} & & \mathbf{0} & & \ddots & 1 \end{pmatrix}. \quad (3.56)$$

Аналогично, обнуление поддиагональных элементов j -го столбца матрицы СЛАУ и соответствующие преобразования правой части можно интерпретировать как умножение системы слева на матрицу

$$E_j = \begin{pmatrix} 1 & & & & & & \mathbf{0} & & \\ & \ddots & & & & & & & \\ & & 1 & & & & & & \\ & & & r_{j+1,j} & 1 & & & & \\ \mathbf{0} & & & \vdots & & \ddots & & & \\ & & & & r_{nj} & & & & 1 \end{pmatrix}. \quad (3.57)$$

В целом метод Гаусса представляется как последовательность умножений обеих частей решаемой СЛАУ слева на матрицы E_j вида (3.57), $j = 1, 2, \dots, n - 1$. При этом матрицей системы становится матрица

$$E_{n-1} \cdots E_2 E_1 A = U, \quad (3.58)$$

которая является верхней треугольной матрицей.

Коль скоро все E_j — нижние треугольные матрицы, их произведение также является нижним треугольным. Кроме того, все E_j неособенны (нижние треугольные с единицами по главной диагонали). Поэтому неособенно и их произведение $E_{n-1} \cdots E_2 E_1$. Вводя обозначение

$$L = (E_{n-1} \cdots E_2 E_1)^{-1},$$

нетрудно понять, что L — нижняя треугольная матрица, для которой в силу (3.58) справедливо

$$A = LU.$$

Получается, что исходная матрица СЛАУ оказалась представленной в виде произведения нижней треугольной L и верхней треугольной U матриц. Это представление называют *треугольным разложением* матрицы или *LU-разложением*.¹²

Соответственно, преобразования матрицы A в прямом ходе метода Гаусса можно трактовать как её разложение на нижний треугольный L и верхний треугольный U множители. В результате исходная СЛАУ представляется в равносильной форме

$$L(Ux) = b,$$

решение которой сводится к решению двух треугольных систем линейных алгебраических уравнений

$$\begin{cases} Ly = b, \\ Ux = y \end{cases} \quad (3.59)$$

с помощью прямой и обратной подстановок соответственно.

Отметим, что при реализации метода Гаусса на компьютере для экономии машинной памяти можно хранить треугольные сомножители L и U на месте A , так как у L по диагонали стоят все единицы.

3.6г Метод Гаусса с выбором ведущего элемента

И в прямом, и в обратном ходе метода Гаусса встречаются операции деления, которые не выполнимы в случае, когда делитель равен нулю. Тогда не может быть выполнен и метод Гаусса в целом. Этот раздел посвящен тому, как модифицировать метод Гаусса, чтобы он был применим для решения любых СЛАУ с неособенными матрицами.

Ведущим элементом в методе Гаусса называют элемент матрицы решаемой системы, на который выполняется деление при исключении

¹²От английских слов lower (нижний) и upper (верхний). Нередко для обозначения этого же понятия можно встретить кальки с иностранных терминов — «LU-факторизация» и «LU-декомпозиция».

поддиагональных элементов очередного столбца.¹³ В алгоритме, описанном в §3.6б, ведущим всюду берётся фиксированный диагональный элемент a_{jj} , вне зависимости от его значения, но желательно модифицировать метод Гаусса так, чтобы ведущий элемент, по возможности, всегда был отличен от нуля. С другой стороны, при решении конкретных СЛАУ, даже в случае $a_{jj} \neq 0$, более предпочтительным иногда может оказаться выбор другого элемента в качестве ведущего.

Отметим, что любое изменение порядка уравнений в системе приводит к равносильной системе уравнений. Но при этом в матрице СЛАУ переставляются строки, так что она существенно меняется. Воспользуемся этим изменением для организации успешного выполнения метода Гаусса.

Назовём *активной подматрицей* j -го шага прямого хода метода Гаусса подматрицу исходной матрицы СЛАУ, образованную последними $n - j + 1$ строками и столбцами. Именно эта подматрица подвергается преобразованиям на j -ом шаге прямого хода, тогда как первые $j - 1$ строк и столбцов остаются уже неизменными.

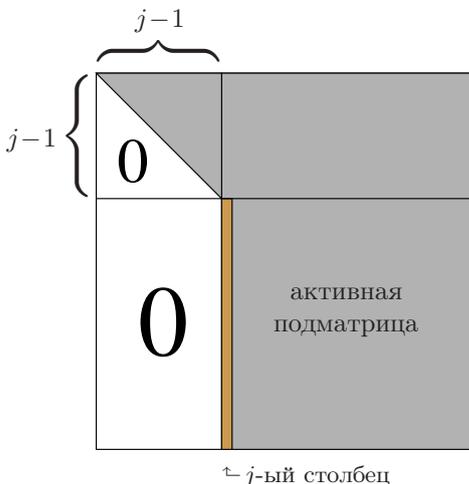


Рис. 3.13. Структура матрицы СЛАУ перед началом j -го шага прямого хода метода Гаусса.

Частичным выбором ведущего элемента на j -ом шаге прямого хо-

¹³Иногда в русской математической литературе его называют *главным* элементом.

да метода Гаусса называют его выбор, как максимального по модулю элемента из всех элементов j -го столбца, лежащих не выше диагонали, и сопровождаемый соответствующей перестановкой строк матрицы и компонент правой части (т. е. уравнений СЛАУ). Максимальным по модулю, а не просто ненулевым, ведущий элемент выбирается для того, чтобы обеспечить наибольшую численную устойчивость алгоритма в условиях вычислений с конечной точностью.

Предложение 3.6.1 *Метод Гаусса с частичным выбором ведущего элемента всегда выполним для систем линейных алгебраических уравнений с неособенными квадратными матрицами.*

Доказательство. Преобразования прямого хода метода Гаусса сохраняют свойство определителя матрицы системы быть неравным нулю. Перед началом j -го шага метода Гаусса эта матрица имеет блочно-треугольный вид, изображённый на Рис. 3.13, и поэтому её определитель равен произведению определителей ведущей подматрицы порядка $(j - 1)$ и активной подматрицы порядка $n - j + 1$. Следовательно, активная подматрица имеет ненулевой определитель, т. е. в первом её столбце обязан найтись хотя бы один ненулевой элемент. Максимальный по модулю из этих ненулевых элементов также ненулевой, и его мы делаем ведущим. Как следствие, прямой ход метода Гаусса выполним.

Обратный ход также не встречает деления на нуль, поскольку полученная в прямом ходе верхняя треугольная матрица неособенна, т. е. все её диагональные элементы должны быть ненулевыми. ■

Каково матричное представление метода Гаусса с выбором ведущего элемента? Введём *элементарные матрицы перестановок*

$$P = \begin{pmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 0 & \cdots & 1 & \\ & & 1 & & & \\ & \vdots & \vdots & \ddots & \vdots & \\ & & & & 1 & \\ 1 & \cdots & & & 0 & \\ & & & & & \ddots \\ & & & & & & 1 \end{pmatrix} \begin{matrix} \\ \\ \leftarrow i\text{-ая строка} \\ \\ \\ \\ \leftarrow j\text{-ая строка} \\ \\ \end{matrix} \tag{3.60}$$

(называемые также *матрицами транспозиции*), которые получаются из единичной матрицы перестановкой двух её строк (или столбцов) с номерами i и j . Умножение на такую матрицу слева приводит к перестановке i -ой и j -ой строк, а умножение справа — к перестановке i -го и j -го столбцов. Тогда для метода Гаусса с частичным выбором ведущего элемента справедливо следующее матричное представление

$$(E_{n-1}P_{n-1}) \cdots (E_1P_1)A = U,$$

где E_j — матрицы преобразований, введённые в предыдущем разделе, а P_1, P_2, \dots, P_{n-1} — элементарные матрицы перестановок, при помощи которых выполняется необходимая перестановка строк на 1-м, 2-м, \dots , $(n-1)$ -м шагах прямого хода метода Гаусса.

Несмотря на то, что метод Гаусса с частичным выбором ведущего элемента теоретически всегда спасает положение, на практике для некоторых «плохих» СЛАУ он может работать не очень хорошо. Это происходит в случае, когда ведущий элемент оказывается малым, и потому коэффициенты r_{ij} из прямого метода Гаусса (3.53) получаются большими по абсолютной величине. По этой причине для обеспечения устойчивости вычислительного процесса по методу Гаусса иногда имеет смысл выбирать ведущий элемент более тщательно, чем это делается при описанном выше частичном выборе.

Заметим, что ещё одним простым способом равносильного преобразования системы уравнений является перенумерация переменных. Ей соответствует перестановка столбцов матрицы, тогда как вектор правых частей при этом неизменен. *Полным выбором* ведущего элемента называют способ его выбора, как максимального по модулю элемента из всей активной подматрицы (а не только из её первого столбца, как было при частичном выборе), и сопровождаемый соответствующей перестановкой строк и столбцов матрицы и компонент правой части. Метод Гаусса с полным выбором ведущего элемента имеет следующее матричное представление

$$(E_{n-1}\check{P}_{n-1}) \cdots (E_1\hat{P}_1)A\hat{P}_1 \cdots \hat{P}_{n-1} = U,$$

где \check{P}_i — элементарные матрицы перестановок, при помощи которых выполняется перестановка строк, \hat{P}_j — элементарные матрицы перестановок, с помощью которых выполняется перестановка столбцов на соответствующих шагах прямого хода метода Гаусса.

Напомним, что *матрицей перестановок* называется матрица, получающаяся из единичной матрицы перестановкой произвольного числа

её строк (или столбцов). Матрица перестановок может быть представлена как произведение нескольких элементарных матриц перестановок вида (3.60) (см. подробности, к примеру, в [7]).

Теорема 3.6.1 *Для неособенной матрицы A существуют матрицы перестановок \check{P} и \hat{P} , такие что*

$$\check{P}A\hat{P} = LU,$$

где L, U — нижняя и верхняя треугольные матрицы, причём диагональными элементами в L являются единицы. В этом представлении можно ограничиться лишь одной из матриц \check{P} или \hat{P} .

Этот результат показывает, что можно один раз переставить строки и столбцы в исходной матрице и потом уже выполнять LU-разложение прямым ходом метода Гаусса без какого-либо специального выбора ведущего элемента. Доказательство теоремы можно найти в [11, 13, 36].

3.6д Существование LU-разложения

В методе Гаусса с выбором ведущего элемента перестановка строк и столбцов приводит к существенному изменению исходной матрицы системы, что не всегда желательно. Естественно задаться вопросом о достаточных условиях реализуемости метода Гаусса без перестановки строк и столбцов. Этот вопрос тесно связан с условиями получения LU-разложения матрицы посредством прямого хода «немодифицированного» метода Гаусса, изложенного в §3.6б.

Теорема 3.6.2 *Если $A = (a_{ij})$ — квадратная $n \times n$ -матрица, у которой все ведущие миноры порядков от 1 до $(n - 1)$ отличны от нуля, т. е.*

$$a_{11} \neq 0, \quad \det \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \neq 0, \quad \dots,$$

$$\det \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1,n-1} \\ a_{21} & a_{22} & \dots & a_{2,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n-1,1} & a_{n-1,2} & \dots & a_{n-1,n-1} \end{pmatrix} \neq 0.$$

то для A существует LU -разложение, т. е. представление её в виде

$$A = LU$$

— произведения нижней треугольной $n \times n$ -матрицы L и верхней треугольной $n \times n$ -матрицы U . Это LU -разложение для A единственно при условии, что диагональными элементами в L являются единицы.

Доказательство проводится индукцией по порядку n матрицы A .

Если $n = 1$, то утверждение теоремы очевидно. Тогда искомые матрицы $L = (l_{ij})$ и $U = (u_{ij})$ являются просто числами, и достаточно взять $l_{11} = 1$ и $u_{11} = a_{11}$.

Пусть теорема верна для матриц размера $(n - 1) \times (n - 1)$. Тогда представим $n \times n$ -матрицу A в блочном виде:

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} = \begin{pmatrix} A_{n-1} & z \\ v & a_{nn} \end{pmatrix},$$

где A_{n-1} — ведущая $(n - 1) \times (n - 1)$ -подматрица A ,

z — вектор-столбец размера $n - 1$,

v — вектор-строка размера $n - 1$,

такие что

$$z = \begin{pmatrix} a_{1n} \\ a_{2n} \\ \vdots \\ a_{n-1,n} \end{pmatrix}, \quad v = (a_{n1} \ a_{n2} \ \dots \ a_{n,n-1}).$$

Требование разложения A на треугольные множители диктует равенство

$$A = \begin{pmatrix} A_{n-1} & z \\ v & a_{nn} \end{pmatrix} = \begin{pmatrix} L_{n-1} & 0 \\ x & l_{nn} \end{pmatrix} \cdot \begin{pmatrix} U_{n-1} & y \\ 0 & u_{nn} \end{pmatrix},$$

где L_{n-1}, U_{n-1} — $(n - 1) \times (n - 1)$ -матрицы,

x — вектор-строка размера $n - 1$,

y — вектор-столбец размера $n - 1$.

Следовательно, используя правила перемножения матриц по блокам, необходимо имеем

$$A_{n-1} = L_{n-1}U_{n-1}, \quad (3.61)$$

$$z = L_{n-1}y, \quad (3.62)$$

$$v = xU_{n-1}, \quad (3.63)$$

$$a_{nn} = xy + l_{nn}u_{nn}. \quad (3.64)$$

Первое из полученных соотношений выполнено в силу индукционного предположения, причём оно должно однозначно определять L_{n-1} и U_{n-1} , если потребовать по диагонали в L_{n-1} единичные элементы. Далее, по условию теоремы $\det A_{n-1} \neq 0$, а потому матрицы L_{n-1} и U_{n-1} также должны быть неособенны. По этой причине системы линейных уравнений относительно x и y —

$$xU_{n-1} = v \quad \text{и} \quad L_{n-1}y = z,$$

которыми являются равенства (3.62)–(3.63), однозначно разрешимы. Стоит отметить, что именно в этом месте доказательства неявно используется условие теоремы, которое требует, чтобы в матрице A все ведущие миноры порядков, меньших чем n , были ненулевыми.

Найдя из (3.62)–(3.63), векторы x и y , мы сможем из соотношения (3.64) восстановить l_{nn} и u_{nn} . Если дополнительно потребовать $l_{nn} = 1$, то значение u_{nn} находится однозначно и равно $(a_{nn} - xy)$. ■

В Теореме 3.6.2 не требуется неособенность всей матрицы A . Из доказательства нетрудно видеть, что при наложенных на A условиях её LU -разложение будет существовать даже при $\det A = 0$, но тогда в матрице U последний элемент u_{nn} будет равен нулю.

В связи с матрицами, имеющими ненулевые ведущие миноры, полезно следующее

Определение 3.6.1 *Квадратная $n \times n$ -матрица $A = (a_{ij})$ называется строго регулярной (или строго неособенной), если все её ведущие миноры, включая и определитель самой матрицы, отличны от нуля, т. е.*

$$a_{11} \neq 0, \quad \det \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \neq 0, \quad \dots, \quad \det A \neq 0.$$

Теорема 3.6.3 Пусть A — квадратная неособенная матрица. Для существования её LU -разложения необходимо и достаточно, чтобы она была строго регулярной.

Доказательство. Достаточность мы доказали в Теореме 3.6.2.

Для доказательства необходимости привлечём блочное представление треугольного разложения $A = LU$. Задавая разные размеры блоков в матрицах A , L и U , получим равенства, аналогичные (3.61). Они означают, что любая ведущая подматрица в A есть произведение ведущих подматриц соответствующих размеров из L и U . Но L и U — неособенные треугольные матрицы, так что все их ведущие подматрицы также неособенны. Отсюда можно заключить неособенность всех ведущих подматриц в A , т. е. её строгую регулярность. ■

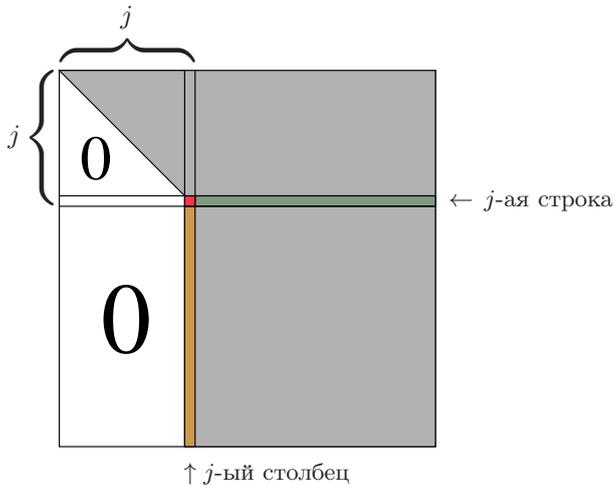


Рис. 3.14. Структура матрицы СЛАУ перед началом j -го шага прямого хода метода Гаусса: другой вид.

В формулировке Теоремы 3.6.2 ничего не говорится о том, реализуем ли метод Гаусса для соответствующей системы линейных алгебраических уравнений. Но нетрудно понять, что в действительности требуемое Теоремой 3.6.2 условие отличия от нуля ведущих миноров в матрице СЛАУ является достаточным для выполнимости рассмотренного в §3.6б варианта метода Гаусса.

Предложение 3.6.2 *Если в системе линейных алгебраических уравнений $Ax = b$ матрица A — квадратная и строго регулярная, то метод Гаусса реализуем в применении к этой системе без перестановки строк и столбцов.*

Доказательство. В самом деле, к началу j -го шага прямого хода, на котором предстоит обнулить поддиагональные элементы j -го столбца матрицы СЛАУ, её ведущей $j \times j$ -подматрицей является треугольная матрица, которая получена из исходной ведущей подматрицы преобразованиями предыдущих $j - 1$ шагов метода Гаусса (см. Рис. 3.14). Эти преобразования — линейное комбинирование строк — не изменяют свойство определителя матрицы быть неравным нулю. Поэтому отличие от нуля какого-либо ведущего минора влечёт отличие от нуля всех диагональных элементов ведущей треугольной подматрицы преобразованной матрицы СЛАУ. В частности, при этом всегда $a_{jj} \neq 0$, так что деление на этот элемент в алгоритмах (3.53) и (3.54) выполнимо. ■

В общем случае проверка как условий Теоремы 3.6.2, так и строгой регулярности матрицы являются весьма непростыми, поскольку вычисление ведущих миноров матрицы требует немалых трудозатрат, и, по существу, ничуть не проще самого метода Гаусса. Тем не менее, условия Теоремы 3.6.2 заведомо выполнены, к примеру, в двух важных частных случаях:

- для СЛАУ с положительно определёнными матрицами (в силу известного критерия Сильвестера),
- если матрица СЛАУ имеет диагональное преобладание (см. признак Адамара, §3.2e).

3.6e Разложение Холецкого

Напомним, что квадратная $n \times n$ -матрица называется *положительно определённой*, если $\langle Ax, x \rangle > 0$ для любых n -векторов x . Ясно, что положительно-определённые матрицы неособенны.

Теорема 3.6.4 *Матрица A является симметричной положительно определённой тогда и только тогда, когда существует неособенная нижняя треугольная матрица C , такая что $A = CC^T$. При этом матрица C из выписанного представления единственна.*

Определение 3.6.2 Представление $A = CC^\top$ называется разложением Холесского, а нижняя треугольная матрица C — множителем Холесского для A .

Доказательство. Пусть $A = CC^\top$ и C неособенна. Тогда неособенна матрица C^\top , и для любого ненулевого вектора $x \in \mathbb{R}^n$ имеем

$$\begin{aligned} \langle Ax, x \rangle &= (Ax)^\top x = (CC^\top x)^\top x \\ &= x^\top CC^\top x = (C^\top x)^\top (C^\top x) = \|C^\top x\|_2^2 > 0, \end{aligned}$$

поскольку $C^\top x \neq 0$. Кроме того, A симметрична по построению. Таким образом, она является симметричной положительно определённой матрицей.¹⁴

Обратно, пусть матрица A симметрична и положительно определена. В силу критерия Сильвестера все её ведущие миноры положительны, а потому на основании Теоремы 3.6.2 о существовании LU-разложения мы можем заключить, что $A = LU$ для некоторых неособенных нижней треугольной матрицы $L = (l_{ij})$ и верхней треугольной матрицы U . Мы дополнительно потребуем, чтобы все диагональные элементы l_{ii} в L были единицами, так что это разложение будет даже однозначно определённым.

Так как

$$LU = A = A^\top = (LU)^\top = U^\top L^\top,$$

то

$$U = L^{-1}U^\top L^\top, \quad (3.65)$$

и далее

$$U(L^\top)^{-1} = L^{-1}U^\top.$$

Слева в этом равенстве стоит произведение верхних треугольных матриц, а справа — произведение нижних треугольных. Равенство, следовательно, возможно лишь в случае, когда левая и правая его части — это диагональная матрица, которую мы обозначим через $D := \text{diag}\{d_1, d_2, \dots, d_n\}$. Тогда из (3.65) вытекает

$$U = L^{-1}U^\top L^\top = DL^\top,$$

¹⁴Это рассуждение никак не использует факт треугольности C и на самом деле обосновывает более общее утверждение: произведение матрицы на её транспонированную является симметричной положительно определённой матрицей.

и потому

$$A = LU = LDL^{\top}. \quad (3.66)$$

Ясно, что в силу неособенности L и U матрица D также неособенна, так что по диагонали у неё стоят ненулевые элементы d_i , $i = 1, 2, \dots, n$. Более того, мы покажем, что все d_i положительны.

Из (3.66) следует, что $D = L^{-1}A(L^{\top})^{-1} = L^{-1}A(L^{-1})^{\top}$. Следовательно, для любого ненулевого вектора x

$$\langle Dx, x \rangle = x^{\top}Dx = x^{\top}L^{-1}A(L^{-1})^{\top}x = ((L^{-1})^{\top}x)^{\top}A((L^{-1})^{\top}x) > 0,$$

так как $(L^{-1})^{\top}x \neq 0$ в силу неособенности матрицы $(L^{-1})^{\top}$. Иными словами, диагональная матрица D положительно определена одновременно с A . Но тогда её диагональные элементы обязаны быть положительными, так как в противном случае, если предположить, что $d_i \leq 0$ для некоторого i , то, беря вектор x равным i -му столбцу единичной матрицы, получим

$$\langle Dx, x \rangle = (Dx)^{\top}x = x^{\top}Dx = d_i \leq 0.$$

Это противоречит положительной определённости матрицы D .

Как следствие, из диагональных элементов матрицы D можно извлекать квадратные корни. Если обозначить получающуюся при этом диагональную матрицу через $\sqrt{D} := \text{diag}\{\sqrt{d_1}, \sqrt{d_2}, \dots, \sqrt{d_n}\}$, то окончательно можем взять $C = L\sqrt{D}$. Это представление для множителя Холецкого, в действительности, единственно, так как по A при сделанных нами предположениях единственным образом определяется нижняя треугольная матрица L , а матричные преобразования, приведшие к формуле (3.66) и её следствиям, обратимы и также дают однозначно определённый результат. ■

3.6ж Метод Холецкого

Основной результат предшествующего раздела мотивирует прямой метод решения систем линейных уравнений, который аналогичен методу (3.59) на основе LU-разложения. Именно, если найдено разложение Холецкого для матрицы A , то решение системы $Ax = b$, равносильной $CC^{\top}x = b$, сводится к решению двух треугольных систем линейных уравнений:

$$\begin{cases} Cy = b, \\ C^{\top}x = y. \end{cases} \quad (3.67)$$

Для решения первой системы применяем прямую подстановку, а для решения второй системы — обратную.

Как конструктивно найти разложение Холецкого? Теорема 3.6.4 носит конструктивный характер и в принципе может служить основой для соответствующего алгоритма. Недостатком этого подхода является существенная опора на LU-разложение матрицы, и потому желательно иметь более прямой способ нахождения разложения Холецкого.

Выпишем равенство $A = CC^T$, определяющее множитель Холецкого, в развёрнутой форме с учётом симметричности A :

$$\begin{pmatrix} a_{11} & & & \\ a_{21} & a_{22} & & \\ \vdots & \vdots & \ddots & \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} = \quad (3.68)$$

$$\begin{pmatrix} c_{11} & & & \\ c_{21} & c_{22} & & \\ \vdots & \vdots & \ddots & \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{pmatrix} \cdot \begin{pmatrix} c_{11} & c_{21} & \dots & c_{n1} \\ & c_{22} & \dots & c_{n2} \\ & & \ddots & \vdots \\ \mathbf{0} & & & c_{nn} \end{pmatrix}, \quad (3.69)$$

где « ∇ » означает симметричные относительно главной диагонали элементы матрицы, которые несущественны в последующих рассуждениях. Можно рассматривать это равенство как систему уравнений относительно неизвестных переменных $c_{11}, c_{21}, c_{22}, \dots, c_{nn}$ — элементов нижнего треугольника множителя Холецкого. Всего их $1+2+\dots+n = \frac{1}{2}n(n+1)$ штук, и для их определения имеем столько же соотношений, вытекающих в этом матричном равенстве из выражений для элементов a_{ij} , $i \geq j$, которые образуют нижний треугольник симметричной матрицы $A = (a_{ij})$.

В поэлементной форме система уравнений (3.68) имеет вид, определяемый правилом умножения матриц и симметричностью A :

$$a_{ij} = \sum_{k=1}^j c_{ik}c_{jk} \quad \text{при } i \geq j. \quad (3.70)$$

Выписанные соотношения образуют, фактически, двумерный массив, в котором уравнения имеют двойные индексы, но их можно линейно упорядочить таким образом, что система уравнений (3.70) получит

специальный вид, очень напоминающий треугольные СЛАУ. Далее эта система может быть решена с помощью процесса, сходного с прямой подстановкой для треугольных СЛАУ (см. §3.6а).

В самом деле, если выписывать выражения для элементов a_{ij} по столбцам матрицы A , начиная в каждом столбце с диагонального элемента a_{jj} и идя сверху вниз до a_{jn} , то все уравнения из (3.70) разбиваются на n следующих групп

$$\begin{cases} c_{j1}^2 + c_{j2}^2 + \dots + c_{j,j-1}^2 + c_{jj}^2 = a_{jj}, \\ c_{i1}c_{j1} + c_{i2}c_{j2} + \dots + c_{ij}c_{jj} = a_{ij}, \quad i = j + 1, \dots, n, \end{cases} \quad (3.71)$$

$$j = 1, 2, \dots, n,$$

В подробной записи

$$\begin{aligned} \text{при } j = 1 & \quad \begin{cases} c_{11}^2 = a_{11}, \\ c_{i1}c_{11} = a_{i1}, \quad i = 2, 3, \dots, n, \end{cases} \\ \text{при } j = 2 & \quad \begin{cases} c_{21}^2 + c_{22}^2 = a_{22}, \\ c_{i1}c_{21} + c_{i2}c_{22} = a_{i2}, \quad i = 3, 4, \dots, n, \end{cases} \\ \text{при } j = 3 & \quad \begin{cases} c_{31}^2 + c_{32}^2 + c_{33}^2 = a_{33}, \\ c_{i1}c_{31} + c_{i2}c_{32} + c_{i3}c_{33} = a_{i3}, \quad i = 4, 5, \dots, n, \end{cases} \\ & \quad \dots \quad \dots \quad \dots \end{aligned}$$

Получается, что в уравнениях (3.71) для j -го столбца множителя Холесского присутствуют все элементы j -го и предшествующих столбцов, но из них реально неизвестными к моменту обработки j -го столбца (т. е. решения j -ой группы уравнений) являются только $(n - j + 1)$ элементов c_{ij} именно j -го столбца, которые к тому же выражаются несложным образом через известные элементы и друг через друга.

В целом выписанная система уравнений (3.71) действительно имеет очень специальный вид, пользуясь которым можно находить элементы c_{ij} матрицы C последовательно друг за другом по столбцам (см.

$$C = \begin{pmatrix} \downarrow & & & & \\ \downarrow & \downarrow & & & \\ \downarrow & \downarrow & \ddots & & \\ \vdots & \vdots & \ddots & & \\ \curvearrowright & \curvearrowright & \dots & \curvearrowright & \times \end{pmatrix} \begin{matrix} \\ \\ \\ \\ 0 \end{matrix}$$

Рис. 3.15. Схема определения элементов треугольного множителя в методе Холецкого.

Рис. 3.15). Более точно,

$$\text{при } j = 1 \quad \begin{cases} c_{11} = \sqrt{a_{11}}, \\ c_{i1} = a_{i1}/c_{11}, \quad i = 2, 3, \dots, n, \end{cases}$$

$$\text{при } j = 2 \quad \begin{cases} c_{22} = \sqrt{a_{22} - c_{21}^2}, \\ c_{i2} = (a_{i2} - c_{i1}c_{21})/c_{22}, \quad i = 3, 4, \dots, n, \end{cases}$$

$$\text{при } j = 3 \quad \begin{cases} c_{33} = \sqrt{a_{33} - c_{31}^2 - c_{32}^2}, \\ c_{i3} = (a_{i3} - c_{i1}c_{31} - c_{i2}c_{32})/c_{33}, \quad i = 4, 5, \dots, n, \end{cases}$$

и так далее для остальных j . Псевдокод этого процесса приведён в Табл. 3.1, где считается, что если верхний предел суммирования превосходит нижний, то сумма «пуста» и суммирование не выполняется.

Если A — симметричная положительно определённая матрица, то в силу Теоремы 3.6.4 система (3.71) обязана иметь решение, и наш алгоритм успешно прорабатывает до конца, находя его. Если же матрица A не является положительно определённой, то алгоритм (3.72) аварийно прекращает работу при попытке извлечь корень из отрицательного числа либо разделить на ноль. Вообще, запуск алгоритма (3.72) — это самый экономичный способ проверки положительной определённости симметричной матрицы.

Метод решения СЛАУ, основанный на разложении Холецкого и использующий соотношения (3.67) и алгоритм (3.72), называют *методом Холецкого*. Он был предложен в 1910 году А.-Л. Холецким в неопубликованной рукописи, которая, тем не менее, сделалась широко извест-

Таблица 3.1. Алгоритм разложения Холецкого

<pre> DO FOR $j = 1$ TO n $c_{jj} \leftarrow \sqrt{a_{jj} - \sum_{k=1}^{j-1} c_{jk}^2}$ DO FOR $i = j + 1$ TO n $c_{ij} \leftarrow \left(a_{ij} - \sum_{k=1}^{j-1} c_{ik}c_{jk} \right) / c_{jj}$ END DO END DO </pre>	(3.72)
---	--------

ной во французской геодезической службе, где решались такие системы уравнений. Позднее метод неоднократно перетоткрывался, и потому иногда в связи с ним используются также термины «метод квадратного корня», «метод квадратных корней» или даже другие имена, данные его позднейшими авторами.

Метод Холецкого можно рассматривать как специальную модификацию метода Гаусса, которая требует вдвое меньше времени и памяти ЭВМ, чем обычный метод Гаусса в общем случае. Замечательным свойством метода Холецкого является также то, что обусловленность множителей Холецкого, вообще говоря, является лучшей, чем у матрицы исходной СЛАУ: она равна корню квадратному из обусловленности исходной матрицы СЛАУ (это следует из самого разложения Холецкого). То есть, в отличие от обычного метода Гаусса, треугольные системы линейных уравнений из (3.67), к решению которых сводится задача, менее чувствительны к ошибкам, чем исходная линейная система. В следующем пункте мы увидим, что подобную ситуацию следует рассматривать как весьма нетипичную.

Если при реализации метода Холецкого использовать комплексную арифметику, то извлечение квадратного корня можно выполнять все-

гда, и потому такая модификация применима к симметричным неособенным матрицам, которые не являются положительно определёнными. При этом множители Холецкого становятся комплексными треугольными матрицами.

Другой популярный способ распространения идеи метода Холецкого на произвольные симметричные матрицы состоит в том, чтобы ограничиться разложением (3.66), которое называется *LDL-разложением матрицы*. Если исходная матрица не является положительно определённой, то диагональными элементами в матрице D могут быть отрицательными. Но LDL-разложение столь же удобно для решения систем линейных алгебраических уравнений, как и рассмотренные ранее треугольные разложения. Детали этих построений читатель может найти, к примеру, в [11, 15, 43, 67].

Отметим также, что существует возможность другой организации вычислений при решении системы уравнений (3.70), когда неизвестные элементы $c_{11}, c_{21}, c_{22}, \dots, c_{nn}$ последовательно находятся по строкам множителя Холецкого, а не по столбцам. Этот алгоритм называется *схемой окаймления* [15], и он по своим свойствам примерно эквивалентен рассмотренному выше алгоритму (3.72).

3.7 Прямые методы на основе ортогональных преобразований

3.7а Число обусловленности и матричные преобразования

Пусть матрица A умножается на матрицу B . Как связано число обусловленности произведения AB с числами обусловленности исходных сомножителей A и B ?

Справедливо

$$\begin{aligned}\|AB\| &\leq \|A\| \|B\|, \\ \|(AB)^{-1}\| &= \|B^{-1}A^{-1}\| \leq \|A^{-1}\| \|B^{-1}\|,\end{aligned}$$

и поэтому

$$\text{cond}(AB) = \|(AB)^{-1}\| \|AB\| \leq \text{cond } A \cdot \text{cond } B. \quad (3.73)$$

С другой стороны, если $C = AB$, то $A = CB^{-1}$, и в силу доказанного неравенства

$$\text{cond}(A) \leq \text{cond}(C) \cdot \text{cond}(B^{-1}) = \text{cond}(AB) \cdot \text{cond}(B),$$

коль скоро $\text{cond}(B^{-1}) = \text{cond}(B)$. Поэтому

$$\text{cond}(AB) \geq \text{cond}(A)/\text{cond}(B).$$

Аналогичным образом из $B = CA^{-1}$ следует

$$\text{cond}(AB) \geq \text{cond}(B)/\text{cond}(A).$$

Объединяя полученные неравенства, в целом получаем оценку

$$\text{cond}(AB) \geq \max \left\{ \frac{\text{cond}(A)}{\text{cond}(B)}, \frac{\text{cond}(B)}{\text{cond}(A)} \right\}. \quad (3.74)$$

Ясно, что её правая часть не меньше 1.

Неравенства (3.73)–(3.74) кажутся грубыми, но они достижимы. В самом деле, пусть A — неособенная симметричная матрица с собственными значениями $\lambda_1, \lambda_2, \dots$ и спектральным числом обусловленности, равным (стр. 266)

$$\text{cond}_2(A) = \frac{\max_i |\lambda_i(A)|}{\min_i |\lambda_i(A)|}.$$

У матрицы A^2 собственные векторы, очевидно, совпадают с собственными векторами матрицы A , а собственные значения равны $\lambda_1^2, \lambda_2^2, \dots$. Как следствие, числом обусловленности матрицы A^2 становится

$$\text{cond}_2(A) = \frac{\max_i (\lambda_i(A))^2}{\min_i (\lambda_i(A))^2} = \frac{\max_i |\lambda_i(A)|^2}{\min_i |\lambda_i(A)|^2} = \left(\frac{\max_i |\lambda_i(A)|}{\min_i |\lambda_i(A)|} \right)^2,$$

и в верхней оценке (3.73) получаем равенство. Совершенно сходным образом можно показать, что для спектрального числа обусловленности оценка (3.73) достигается также на произведениях вида $A^\top A$.

Нижняя оценка (3.74) достигается, к примеру, при $B = A^{-1}$ для чисел обусловленности, порождённых подчинёнными матричными нормами.

Практически наиболее важной является верхняя оценка (3.73), и она показывает, в частности, что при преобразованиях и разложениях матриц число обусловленности может существенно расти. Рассмотрим, к примеру, решение системы линейных алгебраических уравнений

$Ax = b$ методом Гаусса в его матричной интерпретации. Обнуление поддиагональных элементов первого столбца матрицы A — это умножение исходной СЛАУ слева на матрицу E_1 , имеющую вид (3.56), так что мы получаем систему

$$(E_1A)x = E_1b \quad (3.75)$$

с матрицей E_1A , число обусловленности которой оценивается как

$$\text{cond}(E_1A) \leq \text{cond}(E_1) \text{cond}(A).$$

Перестановка строк или столбцов матрицы, выполняемая для поиска ведущего элемента, может незначительно изменить эту оценку в сторону увеличения, так как матрицы перестановок ортогональны и имеют небольшие числа обусловленности. Далее мы обнуляем поддиагональные элементы второго, третьего и т. д. столбцов матрицы системы (3.75), умножая её слева на матрицы E_2, E_3, \dots, E_{n-1} вида (3.57). В результате получаем верхнюю треугольную систему линейных уравнений

$$Ux = y,$$

в которой $U = E_{n-1} \dots E_2 E_1 A$, $y = E_{n-1} \dots E_2 E_1 b$, и число обусловленности матрицы U оценивается сверху как

$$\text{cond}(U) \leq \text{cond}(A) \cdot \text{cond}(E_1) \cdot \text{cond}(E_2) \cdot \dots \cdot \text{cond}(E_{n-1}). \quad (3.76)$$

Если E_j отлична от единичной матрицы, то $\text{cond}(E_j) > 1$, причём несмотря на специальный вид матриц E_j правая и левая части неравенства (3.76) могут отличаться не очень сильно (см. примеры ниже). Как следствие, обусловленность матриц, в которые матрица A исходной СЛАУ преобразуется на промежуточных шагах прямого хода метода Гаусса, а также обусловленность итоговой верхней треугольной матрицы U могут быть существенно хуже, чем у матрицы A .

Пример 3.7.1 Для 2×2 -матрицы (3.10)

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

число обусловленности равно $\text{cond}_2(A) = 14.93$. Выполнение для неё преобразований прямого хода метода Гаусса приводит к матрице

$$\tilde{A} = \begin{pmatrix} 1 & 2 \\ 0 & -2 \end{pmatrix},$$

число обусловленности которой $\text{cond}_2(\tilde{A}) = 4.27$, т. е. уменьшается.

С другой стороны, для матрицы (3.11)

$$B = \begin{pmatrix} 1 & 2 \\ -3 & 4 \end{pmatrix},$$

число обусловленности $\text{cond}_2(B) = 2.62$. Преобразования метода Гаусса превращают её в матрицу

$$\tilde{B} = \begin{pmatrix} 1 & 2 \\ 0 & 10 \end{pmatrix},$$

для которой число обусловленности уже равно $\text{cond}_2(\tilde{B}) = 10.4$, т. е. существенно возрастает.

Числовые данные этого примера читатель может проверить с помощью систем компьютерной математики, таких как Scilab, MATLAB и им аналогичных. ■

Фактически, ухудшение обусловленности и, как следствие, всё большая чувствительность решения к погрешностям в данных — это дополнительная плата за приведение матрицы (и всей СЛАУ) к удобному для решения виду. Можно ли уменьшить эту плату? И если да, то как?

Хорошей идеей является привлечение для матричных преобразований ортогональных матриц, которые имеют наименьшую возможную обусловленность в спектральной норме (и небольшие числа обусловленности в других нормах). Умножение на такие матрицы, по крайней мере, не будет ухудшать обусловленность получающихся систем линейных уравнений и устойчивость их решений к ошибкам вычислений.

3.7б QR-разложение матриц

Определение 3.7.1 Для матрицы A представление $A = QR$ в виде произведения ортогональной матрицы Q и правой треугольной матрицы R называется QR-разложением.

По поводу этого определения следует пояснить, что правая треугольная матрица — это то же самое, что верхняя треугольная матрица, которую мы условились обозначать U . Другая терминология обусловлена здесь историческими причинами, и частичное её оправдание состоит в том, что QR-разложение матрицы действительно «совсем

другое», нежели LU-разложение. Впрочем, в математической литературе можно встретить тексты, где LU-разложение матрицы называется «LR-разложением» (от английских слов left-right), т. е. разложением на «левую и правую треугольные матрицы».

Теорема 3.7.1 *QR-разложение существует для любой квадратной матрицы.*

Доказательство. Если A — неособенная матрица, то, как было показано при доказательстве Теоремы 3.6.4, $A^T A$ — симметричная положительно определённая матрица. Следовательно, существует её разложение Холесского

$$A^T A = R^T R,$$

где R — правая (верхняя) треугольная матрица. При этом R , очевидно, неособенна. Тогда матрица $Q := AR^{-1}$ ортогональна, поскольку

$$\begin{aligned} Q^T Q &= (AR^{-1})^T AR^{-1} = (R^{-1})^T A^T A R^{-1} \\ &= (R^{-1})^T (R^T R) R^{-1} = ((R^{-1})^T R^T)(RR^{-1}) = I. \end{aligned}$$

Следовательно, в целом $A = QR$, где определённые выше сомножители Q и R удовлетворяют условиям теоремы.

Рассмотрим теперь случай особенной матрицы A . Известно, что любую особенную матрицу можно приблизить последовательностью неособенных. Например, это можно сделать с помощью матриц $A_k = A + \frac{1}{k}I$, начиная с достаточно больших натуральных номеров k . При этом собственные значения A_k суть $\lambda(A_k) = \lambda(A) + \frac{1}{k}$, и если величина $\frac{1}{k}$ меньше расстояния от нуля до ближайшего ненулевого собственного значения матрицы A , то A_k неособенна.

В силу уже доказанного для всех матриц из последовательности $\{A_k\}$ существуют QR-разложения:

$$A_k = Q_k R_k,$$

где все Q_k ортогональны, а R_k — правые треугольные матрицы. В качестве ортогонального разложения для A можно было бы взять пределы матриц Q_k и R_k , если таковые существуют. Но сходятся ли куда-нибудь последовательности этих матриц при $k \rightarrow \infty$, когда $A_k \rightarrow A$? Ответ на это вопрос может быть отрицательным, а потому приходится

действовать более тонко, выделяя из $\{A_k\}$ подходящую подпоследовательность.

Множество ортогональных матриц компактно, поскольку является замкнутым (прообраз единичной матрицы I при непрерывном отображении $X \mapsto X^T X$) и ограничено ($\|X\|_2 \leq 1$). Поэтому из последовательности ортогональных матриц $\{Q_k\}$ можно выбрать сходящуюся подпоследовательность $\{Q_{k_l}\}_{l=1}^\infty$. Ей соответствуют подпоследовательности $\{A_{k_l}\}$ и $\{R_{k_l}\}$, причём первая из них также сходится, как подпоследовательность сходящейся последовательности $\{A_k\}$.

Обозначим $Q := \lim_{l \rightarrow \infty} Q_{k_l}$, и это также ортогональная матрица. Тогда

$$\lim_{l \rightarrow \infty} (Q_{k_l}^T A_{k_l}) = \lim_{l \rightarrow \infty} Q_{k_l}^T \cdot \lim_{l \rightarrow \infty} A_{k_l} = Q^T A = R$$

— правой треугольной матрице, поскольку все $Q_{k_l}^T A_{k_l}$ были правыми треугольными матрицами R_{k_l} . Таким образом, в целом снова $A = QR$ с ортогональной Q и правой треугольной R , как и требовалось. ■

Если известно QR-разложение матрицы A , то решение исходной СЛАУ, равносильной

$$(QR)x = b$$

сводится к решению треугольной системы линейных алгебраических уравнений

$$Rx = Q^T b. \quad (3.77)$$

Ниже в §3.17e мы встретимся и с другими важными применениями QR-разложения матриц — при численном решении проблемы собственных значений.

Хотя для неособенных матриц доказательство Теоремы 3.7.1 носит конструктивный характер, оно существенно завязано на разложение Холецкого матрицы $A^T A$, а потому находить с его помощью QR-разложение не очень удобно. На практике основным инструментом получения QR-разложения является техника, использующая так называемые матрицы отражения и матрицы вращения, описанию которых посвящены следующие разделы книги.

3.7в Ортогональные матрицы отражения

Определение 3.7.2 Для вектора $u \in \mathbb{R}^n$ с единичной евклидовой нормой, $\|u\|_2 = 1$, матрица $H = H(u) = I - 2uu^T$ называется матрицей

отражения или матрицей Хаусхолдера. Вектор u называется порождающим или вектором Хаусхолдера для матрицы отражения $H(u)$.

Предложение 3.7.1 Матрицы отражения являются симметричными ортогональными матрицами. Кроме того, для матрицы $H(u)$

порождающий вектор u является собственным вектором, отвечающим собственному значению (-1) , т. е. $H(u) \cdot u = -u$;

любой вектор $v \in \mathbb{R}^n$, ортогональный порождающему вектору u , является собственным вектором, отвечающим собственному значению 1 , т. е. $H(u) \cdot v = v$.

Доказательство проводится непосредственной проверкой.

Симметричность матрицы $H(u)$:

$$\begin{aligned} H^\top &= (I - 2uu^\top)^\top = I^\top - (2uu^\top)^\top \\ &= I - 2(u^\top)^\top u^\top = I - 2uu^\top = H. \end{aligned}$$

Ортогональность:

$$\begin{aligned} H^\top H &= (I - 2uu^\top)(I - 2uu^\top) \\ &= I - 2uu^\top - 2uu^\top + 4uu^\top uu^\top \\ &= I - 4uu^\top + 4u(u^\top u)u^\top = I, \quad \text{так как } u^\top u = 1. \end{aligned}$$

Собственные векторы и собственные значения:

$$H(u) \cdot u = (I - 2uu^\top)u = u - 2u(u^\top u) = u - 2u = -u;$$

$$H(u) \cdot v = (I - 2uu^\top)v = v - 2u(u^\top v) = v, \quad \text{поскольку } u^\top v = 0.$$

Это завершает доказательство предложения. ■

Из последних двух свойств матриц отражения следует геометрическая интерпретация, которая мотивирует их название. Эти матрицы действительно осуществляют преобразование отражения относительно гиперплоскости, ортогональной порождающему вектору u .

Чтобы убедиться в этом, представим произвольный вектор x в виде $\alpha u + v$, где u — порождающий матрицу отражения вектор, а v — ему ортогональный, т. е. $u^\top v = 0$ (см. Рис. 3.16). Тогда

$$H(u) \cdot x = H(u) \cdot (\alpha u + v) = -\alpha u + v,$$

т.е. в преобразованном матрицей $H(u)$ векторе компонента, ортогональная рассматриваемой гиперплоскости, сменила направление на противоположное. Это и соответствует отражению относительно неё.

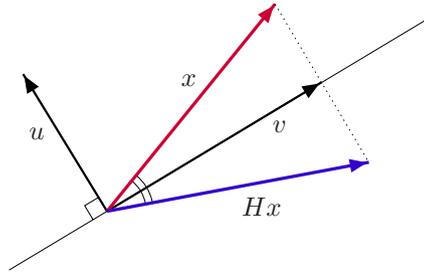


Рис. 3.16. Геометрическая интерпретация действия матрицы отражения.

Предложение 3.7.2 Для любого ненулевого вектора $x \in \mathbb{R}^n$ существует матрица отражения, переводящая его в вектор, коллинеарный заданному вектору $e \in \mathbb{R}^n$ с единичной длиной, $\|e\|_2 = 1$.

Доказательство. Если H — искомая матрица отражения, и u — порождающий её вектор Хаусхолдера, то утверждение предложения требует равенства

$$Hx = x - 2(uu^\top)x = \gamma e \quad (3.78)$$

с некоторым коэффициентом $\gamma \neq 0$. Отдельно рассмотрим два случая — когда векторы x и e неколлинеарны, и когда они коллинеарны друг другу.

В первом случае можно переписать (3.78) в виде равенства

$$2u(u^\top x) = x - \gamma e, \quad (3.79)$$

правая часть которого заведомо не равна нулю. Тогда и числовой множитель $u^\top x$ в левой части обязан быть ненулевым, и из соотношения (3.79) можно заключить, что

$$u = \frac{1}{2u^\top x} (x - \gamma e),$$

т.е. что вектор u , порождающий искомую матрицу отражения, должен быть коллинеарен вектору $(x - \gamma e)$.

Для определения коэффициента γ заметим, что ортогональная матрица H не изменяет длин векторов, так что $\|Hx\|_2 = \|x\|_2$. С другой стороны, взяв евклидову норму от обеих частей (3.78), получим $\|Hx\|_2 = |\gamma| \|e\|_2$. Сопоставляя оба равенства, можем заключить

$$\|x\|_2 = |\gamma| \|e\|_2, \quad \text{т. е. } \gamma = \pm \|x\|_2.$$

Следовательно, вектор Хаусхолдера u коллинеарен векторам

$$\tilde{u} = x \pm \|x\|_2 e, \quad (3.80)$$

и для окончательного определения u остаётся лишь применить нормировку:

$$u = \frac{\tilde{u}}{\|\tilde{u}\|_2}.$$

Тогда $H = I - 2uu^\top$ — искомая матрица отражения.

Обсудим теперь случай, когда x коллинеарен e . При этом предшествующая конструкция частично теряет смысл, так как вектор $\tilde{u} = x - \gamma e$ может занулиться при подходящем выборе множителя γ .

Но даже если $x - \gamma e = 0$ для какого-то одного из значений $\gamma = -\|x\|_2$ и $\gamma = \|x\|_2$, то для противоположного по знаку значения γ наверняка $x - \gamma e \neq 0$. Более формально можно сказать, что конкретный знак у множителя $\gamma = \pm \|x\|_2$ следует выбирать из условия максимизации нормы вектора $(x - \gamma e)$. Далее все рассуждения, следующие за формулой (3.79), остаются в силе и приводят к определению вектора Хаусхолдера.

Наконец, в случае коллинеарных векторов x и e мы можем просто указать явную формулу для вектора Хаусхолдера:

$$u = \frac{x}{\|x\|_2}.$$

При этом

$$u^\top x = \frac{x^\top x}{\|x\|_2} = \|x\|_2 \neq 0,$$

и для соответствующей матрицы отражения имеет место

$$Hx = x - 2(uu^\top)x = x - 2u(u^\top x) = x - 2 \frac{x}{\|x\|_2} \|x\|_2 = -x.$$

Итак, вектор x снова переводится матрицей H в вектор, коллинеарный вектору e , т. е. условие предложения удовлетворено и в этом случае.¹⁵

■

В доказательстве предложения присутствовала неоднозначность в выборе знака в выражении $\tilde{u} = x \pm \|x\|_2 e$, если x и e неколлинеарны. В действительности, годится любой знак, и его конкретный выбор может определяться, как мы увидим, требованием устойчивости вычислительного алгоритма.

3.7г Метод Хаусхолдера

В основе *метода Хаусхолдера* для решения систем линейных алгебраических уравнений (который называют также *методом отражений*) лежит та же самая идея, что и в методе Гаусса: привести эквивалентными преобразованиями исходную систему к правому (верхнему) треугольному виду, а затем воспользоваться обратной подстановкой (3.54). Но теперь это приведение выполняется более глубокими, чем в методе Гаусса, преобразованиями матрицы, именно, путём последовательного умножения на специальным образом подобранные матрицы отражения.

Предложение 3.7.3 *Для любой квадратной матрицы A существует конечная последовательность H_1, H_2, \dots, H_{n-1} , состоящая из матриц отражения и, возможно, единичных матриц, таких что матрица*

$$H_{n-1}H_{n-2} \cdots H_2H_1A = R$$

является правой треугольной матрицей.

Раздельное упоминание матриц отражения и единичных матриц вызвано здесь тем, что единичная матрица не является матрицей отражения.

Для формального описания алгоритма очень удобно применять систему обозначений матрично-векторных объектов, укоренившуюся в языках программирования высокого уровня Fortran, MATLAB, Scilab и

¹⁵Интересно, что этот тонкий случай доказательства имеет, скорее, теоретическое значение, так как на практике если вектор уже имеет нужное направление, то с ним, как правило, можно вообще ничего не делать.

др. В частности, посредством $A(p : q, r : s)$ обозначается *сечение* массива A , которое определяется как массив с тем же количеством измерений и имеющий элементы, которые стоят на пересечении строк с номерами с p по q и столбцов с номерами с r по s . То есть, запись $A(p : q, r : s)$ указывает в индексах матрицы A не отдельные значения, а целые диапазоны изменения индексов элементов, из которых образуется новая матрица, как подматрица исходной.

Доказательство предложения конструктивно.

Используя результат Предложения 3.7.2, возьмём в качестве H_1 матрицу отражения, которая переводит 1-й столбец A в вектор, коллинеарный $(1, 0, \dots, 0)^T$, если хотя бы один из элементов $a_{21}, a_{31}, \dots, a_{n1}$ не равен нулю. Иначе полагаем $H_1 = I$. Затем переходим к следующему шагу.

В результате выполнения первого шага матрица СЛАУ приводится, как и в методе Гаусса, к виду

$$\tilde{A} = \left(\begin{array}{c|cccc} \times & \times & \times & \cdots & \times \\ \hline 0 & \times & \times & \cdots & \times \\ 0 & \times & \times & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \times & \times & \cdots & \times \end{array} \right).$$

где крестиками « \times » обозначены элементы, которые, возможно, не равны нулю. Прделаем теперь то же самое с матрицей $\tilde{A}(2 : n, 2 : n)$, обнулив у неё подходящим отражением поддиагональные элементы первого столбца, который является вторым во всей большой матрице. И так далее до $(n - 1)$ -го столбца.

Для формального описания алгоритма определим теперь матрицу $H_j = H_j(u)$, $j = 2, 3, \dots, n - 1$, как $n \times n$ -матрицу отражения, порождаемую вектором Хаусхолдера $u \in \mathbb{R}^n$, который имеет нулевыми первые $j - 1$ компонент и подобран так, чтобы $H_j(u)$ аннулировала поддиагональные элементы j -го столбца в матрице $\tilde{A} = H_{j-1} \cdots H_2 H_1 A$, если среди этих поддиагональных элементов существуют ненулевые. Иначе, если в преобразуемой матрице $\tilde{A} = (\tilde{a}_{ij})$ все элементы $\tilde{a}_{j+1,j}, \tilde{a}_{j+2,j}, \dots, \tilde{a}_{nj}$ — нулевые, то полагаем $H_j = I$ — единичной $n \times n$ -матрице.

Таблица 3.2. QR-разложение матрицы
с помощью отражений Хаусхолдера

```

DO FOR  $j = 1$  TO  $n - 1$ 
  IF ( вектор  $A((j + 1) : n, j)$  ненулевой ) THEN
    вычислить вектор Хаусхолдера  $u \in \mathbb{R}^{n-j+1}$ ,
    порождающий отражение, которое переводит
    вектор  $A(j : n, j)$  в вектор, коллинеарный
    вектору  $(1, 0, \dots, 0)^\top$ ;
     $\tilde{H} \leftarrow I - 2uu^\top$ ;
  ELSE
     $\tilde{H} \leftarrow I$ ;
  END IF
   $A(j : n, j : n) \leftarrow \tilde{H} A(j : n, j : n)$ ;
END DO

```

Можно положить в блочной форме

$$H_i = \left(\begin{array}{c|c} I & 0 \\ \hline 0 & \tilde{H}_i \end{array} \right),$$

где в верхнем левом углу стоит единичная $(j - 1) \times (j - 1)$ -матрица, а \tilde{H}_i — матрица размера $(n - j + 1) \times (n - j + 1)$, которая переводит вектор $\tilde{A}(j : n, j)$ в $(n - j + 1)$ -вектор, коллинеарный $(1, 0, \dots, 0)^\top$, т. е. обнуляет поддиагональные элементы j -го столбца в \tilde{A} . Если хотя бы один из элементов $\tilde{a}_{j+1,j}, \tilde{a}_{j+2,j}, \dots, \tilde{a}_{nj}$ не равен нулю, то \tilde{H}_i — матрица отражения, способ построения которой описывается в Предложении 3.7.2. Иначе, если $(\tilde{a}_{j+1,j}, \tilde{a}_{j+2,j}, \dots, \tilde{a}_{nj})^\top = 0$, то \tilde{H}_i единичная $(n - j + 1) \times (n - j + 1)$ -матрица. ■

Отметим, что из представления

$$H_{n-1}H_{n-2} \cdots H_2H_1A = R$$

вытекает равенство $A = QR$ с ортогональной матрицей

$$Q = (H_{n-1}H_{n-2} \cdots H_2H_1)^{-1}.$$

Таким образом, мы получаем QR-разложение матрицы A , т.е. Предложения 3.7.2 и 3.7.3 дают в совокупности ещё одно, конструктивное, доказательство Теоремы 3.7.1. Соответствующий псевдокод алгоритма для вычисления QR-разложения матрицы приведён в Табл. 3.2.

Как следствие, исходная система уравнений $Ax = b$ становится равносильной системе уравнений

$$\begin{cases} Qy = b, \\ Rx = y, \end{cases}$$

с несложно решаемыми составными частями. При практической реализации удобнее дополнить алгоритм Табл. 3.2 инструкциями, которые задают преобразования вектора правой части СЛАУ, и тогда результатом работы нового алгоритма будет правая треугольная СЛАУ $Rx = y$. Её можно решать с помощью обратной подстановки (3.54).

Согласно Предложению 3.7.2 вычисление вектора Хаусхолдера u в качестве первого шага требует нахождения из (3.80) вектора \tilde{u} , в котором имеется неоднозначность выбора знака второго слагаемого. При вычислениях на цифровых ЭВМ в стандартной арифметике с плавающей точкой имеет смысл брать

$$\tilde{u} = \begin{cases} A(j : n, j) + \|A(j : n, j)\|_2 e, & \text{если } a_{jj} \geq 0, \\ A(j : n, j) - \|A(j : n, j)\|_2 e, & \text{если } a_{jj} \leq 0, \end{cases}$$

где $e = (1, 0, \dots, 0)^\top$. Тогда вычисление первого элемента в столбце $A(j : n, j)$, т.е. того единственного элемента, который останется ненулевым, не будет сопровождаться вычитанием чисел одного знака и, как следствие, возможной потерей точности.

Ещё одно соображение по практической реализации описанного в Предложении 3.7.3 алгоритма состоит в том, что в действительности даже не нужно формировать в явном виде матрицу отражения \tilde{H} : умножение на неё можно выполнить по экономичной формуле

$$\begin{aligned} (I - 2uu^\top) A(j : n, j : n) \\ = A(j : n, j : n) - 2u (u^\top A(j : n, j : n)). \end{aligned}$$

Определённым недостатком метода Хаусхолдера и описываемого в следующем пункте метода вращений в сравнении с методом Гаусса является привлечение неарифметической операции извлечения квадратного корня, которая приводит к иррациональностям. Это не позволяет точно (без округлений) реализовать соответствующие алгоритмы в поле рациональных чисел, к примеру, в программных системах так называемых «безошибочных вычислений» или языках программирования типа Ruby [82], которые могут оперировать рациональными дробями с числителем и знаменателем в виде целых чисел.

3.7д Матрицы вращения

Пусть даны натуральные числа k, l , не превосходящие n , т. е. размерности пространства \mathbb{R}^n , и задано значение угла θ , $0 \leq \theta < 2\pi$. Матрицей вращения называется $n \times n$ -матрица $G(k, l, \theta)$ вида

$$\begin{array}{l}
 k\text{-ая строка} \\
 \\
 \\
 \\
 \\
 \\
 l\text{-ая строка}
 \end{array}
 \left(
 \begin{array}{cccc}
 1 & & & \\
 & \ddots & & \\
 & & \cos \theta & \cdots & -\sin \theta \\
 & & & 1 & \\
 & & \vdots & \ddots & \vdots \\
 & & & & 1 \\
 & & \sin \theta & \cdots & \cos \theta \\
 & & & & & \ddots \\
 & & & & & & 1
 \end{array}
 \right), \quad (3.81)$$

где все не выписанные явно элементы вне главной диагонали равны нулю. Таким образом, $G(k, l, \theta)$ — это матрица, которая отличается от единичной матрицы лишь элементами, находящимися в позициях (k, k) , (k, l) , (l, k) и (l, l) для данных индексов k, l . Нетрудно проверить, что она ортогональна.

Матрица

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

задаёт, как известно, вращение двумерной плоскости $0x_1x_2$ на угол θ вокруг начала координат.¹⁶ Матрица $G(k, l, \theta)$ также задаёт вращение

¹⁶Напомним, что положительным направлением вращения плоскости считается вращение «против часовой стрелки».

пространства \mathbb{R}^n на угол θ вокруг оси, проходящей через начало координат и ортогональной гиперплоскости $0x_kx_l$. Матрицы вращения $G(k, l, \theta)$ называют также *матрицами Гивенса*, и мы будем иногда обозначать их посредством $G(k, l)$, если конкретная величина угла θ несущественна.

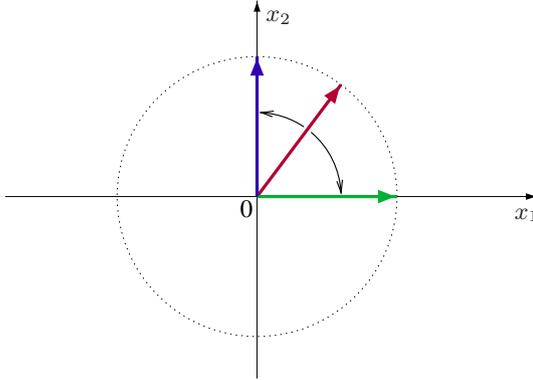


Рис. 3.17. Подходящим вращением можно занулить любую из компонент двумерного вектора.

Если вектор $a = (a_1, a_2)^\top$ — ненулевой, то, взяв

$$\cos \theta = \frac{a_1}{\|a\|_2}, \quad \sin \theta = \frac{-a_2}{\|a\|_2}, \quad \text{где } \|a\|_2 = \sqrt{a_1^2 + a_2^2},$$

мы можем с помощью матрицы двумерного вращения занулить вторую компоненту этого вектора:

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} \|a\|_2 \\ 0 \end{pmatrix}.$$

Аналогично может быть занулена первая компонента вектора a , путём домножения на такую матрицу вращения, что

$$\cos \theta = \frac{a_2}{\|a\|_2}, \quad \sin \theta = \frac{a_1}{\|a\|_2}.$$

В общем случае умножение любой матрицы $A = (a_{ij})$ слева на матрицу вращения $G(k, l, \theta)$ приводит к тому, что в их произведении

$\tilde{A} = (\tilde{a}_{ij}) := G(k, l, \theta) A$ строки k -ая и l -ая становятся линейными комбинациями строк с этими же номерами из A :

$$\begin{aligned}\tilde{a}_{kj} &\leftarrow a_{kj} \cos \theta - a_{lj} \sin \theta, \\ \tilde{a}_{lj} &\leftarrow a_{kj} \sin \theta + a_{lj} \cos \theta,\end{aligned}\tag{3.82}$$

$j = 1, 2, \dots, n$. Остальные элементы матрицы \tilde{A} совпадают с элементами матрицы A . Из рассуждений предшествующего абзаца вытекает, что путём специального подбора угла θ можно всегда занулить элемент в произвольной наперёд заданной позиции k -ой или l -ой строки матрицы $\tilde{A} = G(k, l, \theta) A$.

Как следствие, любая квадратная матрица A может быть приведена к правому треугольному виду с помощью последовательности умножений слева на матрицы вращения. Более точно, мы можем один за другим занулить поддиагональные элементы первого столбца, потом второго, третьего и т. д., аналогично тому, как это делалось в методе Гаусса. При этом зануление поддиагональных элементов второго и последующих столбцов никак не испортит полученные ранее нулевые элементы предшествующих столбцов, так как линейное комбинирование нулей даст снова нуль. В целом, существует набор матриц вращения $G(1, 2), G(1, 3), \dots, G(1, n), G(2, 3), \dots, G(n-1, n)$, таких что

$$G(n-1, n) \cdots G(2, 3) G(1, n) \cdots G(1, 3) G(1, 2) A = R$$

— правой треугольной матрице. Отсюда

$$A = G(1, 2)^\top G(1, 3)^\top \cdots G(1, n)^\top G(2, 3)^\top \cdots G(n-1, n)^\top R,$$

и мы получили QR-разложение матрицы A , так как произведение транспонированных матриц вращения также является ортогональной матрицей.

Использование преобразований вращения — ещё один конструктивный способ получения QR-разложения, технически даже более простой, чем метод отражений Хаусхолдера. При его реализации организовывать полноценные матрицы вращения $G(k, l, \theta)$ и матричные умножения с ними, конечно, нецелесообразно, так как большинством элементов в $G(k, l, \theta)$ являются нули. Результат умножения слева на матрицу вращения разумно находить путём перевычисления лишь ненулевых элементов всего двух строк по формулам (3.82).

Для плотно заполненных матриц использование вращений в полтора раза более трудоёмко, чем получение QR-разложения с помощью

матриц отражения, но зато вращения более предпочтительны для разреженных матриц в силу своей большей гибкости при занулении отдельных элементов.

3.7е Процессы ортогонализации

Ортогонализацией называют процесс построения по заданному базису линейного пространства некоторого ортогонального базиса, который имеет ту же самую линейную оболочку. Ввиду удобства ортогональных базисов для представления решений разнообразных задач и, как следствие, их важности во многих приложениях (см., к примеру, §2.10г) огромное значение имеют и процессы ортогонализации.

Исторически первым процессом ортогонализации был алгоритм, который по традиции связывают с именами Й. Грама и Э. Шмидта.¹⁷ По конечной линейно независимой системе векторов v_1, v_2, \dots, v_n процесс Грама-Шмидта строит ортогональный базис q_1, q_2, \dots, q_n линейной оболочки векторов v_1, v_2, \dots, v_n .

Возьмём в качестве первого вектора q_1 конструируемого ортогонального базиса вектор v_1 , первый из исходного базиса. Далее для построения q_2 можно использовать v_2 «как основу», но откорректировав его с учётом требования ортогональности к q_1 и принадлежности линейной оболочке векторов $q_1 = v_1$ и v_2 . Естественно положить $q_2 = v_2 + \alpha_{21}q_1$, где коэффициент α_{21} подлежит определению из условия ортогональности

$$\langle q_1, v_2 + \alpha_{21}q_1 \rangle = 0.$$

Отсюда

$$\alpha_{21} = -\frac{\langle q_1, v_2 \rangle}{\langle q_1, q_1 \rangle}.$$

Далее аналогичным образом находится $q_3 = v_3 + \alpha_{31}q_1 + \alpha_{32}q_2$, и т. д.

В целом ортогонализация Грама-Шмидта выполняется в соответствии со следующими расчётными формулами:

$$q_j \leftarrow v_j - \sum_{k=1}^{j-1} \frac{\langle q_k, v_j \rangle}{\langle q_k, q_k \rangle} q_k, \quad j = 1, 2, \dots, n. \quad (3.83)$$

В Табл. 3.3 дан псевдокод ортогонализации Грама-Шмидта, дополненной ещё нормализаций получающихся векторов.

¹⁷Иногда этот процесс называют «ортогонализацией Сонина-Шмидта».

Таблица 3.3. Ортогонализация Грама-Шмидта

```

DO FOR  $j = 1$  TO  $n$ 
   $q_j \leftarrow v_j$ ;
  DO  $k = 1$  TO  $j - 1$ 
     $\alpha_{kj} \leftarrow \langle q_k, v_j \rangle$ ;
     $q_j \leftarrow q_j - \alpha_{kj} q_k$ ;
  END DO
   $\alpha_{jj} \leftarrow \|q_j\|_2$ ;
  IF ( $\alpha_{jj} = 0$ ) THEN
    STOP, сигнализируя « $v_j$  линейно зависит
      от векторов  $v_1, v_2, \dots, v_{j-1}$ »
  END IF
   $q_j \leftarrow q_j / \alpha_{jj}$ ;
END DO

```

Дадим матричное представление процесса ортогонализации Грама-Шмидта.

Пусть векторы v_1, v_2, \dots, v_n заданы своими координатными представлениями в некотором базисе, и из вектор-столбцов этих координатных представлений мы организуем матрицу W . В результате ортогонализации мы должны получить ортогональную матрицу, в которой первый столбец — это нормированный первый вектор, второй столбец — это нормированная линейная комбинация первых двух вектор-столбцов, и т. д. Столбец с номером j результирующей ортогональной матрицы равен нормированной линейной комбинации первых j штук столбцов исходной матрицы. В целом процесс ортогонализации Грама-Шмидта равносильен умножению W слева на верхнюю треугольную матрицу, в результате чего должна получиться ортогональная матрица.

Фактически, ортогонализацию Грама-Шмидта можно рассматривать как ещё один способ получения QR-разложения матрицы. Но свойства этого процесса существенно хуже, чем у метода отражений или метода вращений. Если исходная система векторов близка к линейно

зависимой, то полученный в результате применения алгоритма Грама-Шмидта базис может существенно отличаться от ортогонального в том смысле, что попарные скалярные произведения его векторов будут заметно отличаться от нуля.

Таблица 3.4. Модифицированный алгоритм ортогонализации Грама-Шмидта

```

DO FOR  $j = 1$  TO  $n$ 
   $q_j \leftarrow v_j$ ;
  DO  $k = 1$  TO  $j - 1$ 
     $\alpha_{kj} \leftarrow \langle q_k, q_j \rangle$ ;
     $q_j \leftarrow q_j - \alpha_{kj} q_k$ ;
  END DO
   $\alpha_{jj} \leftarrow \|q_j\|_2$ ;
  IF ( $\alpha_{jj} = 0$ ) THEN
    STOP, сигнализируя « $v_j$  линейно зависит
    от векторов  $v_1, v_2, \dots, v_{j-1}$ »
  END IF
   $q_j \leftarrow q_j / \alpha_{jj}$ ;
END DO

```

Этот недостаток можно до некоторой степени исправить, модифицировав расчётные формулы алгоритма Грама-Шмидта так, чтобы вычисление поправочных коэффициентов α_{kj} выполнялось другим способом. Псевдокод модифицированной ортогонализации Грама-Шмидта дан в Табл. 3.4.

В общем случае при ортогонализации Грама-Шмидта построение каждого следующего вектора требует привлечения всех ранее построенных векторов. Но если исходная система векторов имеет специальный вид, в определённом смысле согласованный с используемым скалярным произведением, то ситуация упрощается. Важнейший частный случай — ортогонализация так называемых *подпространств Крылова*.

Определение 3.7.3 Пусть A — квадратная $n \times n$ -матрица, r — n -

вектор. Подпространствами Крылова $\mathcal{K}_i(A, r)$, $i = 1, 2, \dots, n$, матрицы A относительно вектора r называются линейные оболочки векторов $r, Ar, \dots, A^{i-1}r$, т. е. $\mathcal{K}_i(A, r) = \text{lin}\{r, Ar, \dots, A^{i-1}r\}$.

Оказывается, что если A — симметричная положительно определённая матрица, то при ортогонализации подпространств Крылова построение каждого последующего вектора привлекает лишь два предшествующих вектора из строящегося базиса. Более точно, справедлива

Теорема 3.7.2 Пусть векторы $r, Ar, A^2r, \dots, A^{n-1}r$ линейно независимы. Если векторы p_0, p_1, \dots, p_{n-1} получены из них с помощью процесса ортогонализации, то они выражаются трёхчленными рекуррентными соотношениями

$$\begin{aligned} p_0 &\leftarrow r, \\ p_1 &\leftarrow Ap_1 - \alpha_1 p_1, \\ p_{k+1} &\leftarrow Ap_k - \alpha_k p_k - \beta_k p_{k-1}, \quad k = 1, 2, \dots, n-2, \end{aligned}$$

где коэффициенты ортогонализации α_k и β_k вычисляются следующим образом:

$$\begin{aligned} \alpha_k &= \frac{\langle Ap_k, p_k \rangle}{\langle p_k, p_k \rangle}, \quad k = 0, 1, \dots, n-2, \\ \beta_k &= \frac{\langle Ap_k, p_{k-1} \rangle}{\langle p_{k-1}, p_{k-1} \rangle} = \frac{\langle p_k, p_k \rangle}{\langle p_{k-1}, p_{k-1} \rangle}, \quad k = 1, 2, \dots, n-2. \end{aligned}$$

Этот факт был открыт К. Ланцошем в 1952 году и имеет многочисленные применения в практике вычислений. В частности, он существенно используется в методе сопряжённых градиентов для решения СЛАУ (см. §3.10г).

Доказательство. Если векторы p_0, p_1, \dots, p_{n-1} получены из $r, Ar, A^2r, \dots, A^{n-1}r$ в результате ортогонализации, то из формул (3.83) следует

$$p_k = A^k r + \sum_{i=0}^{k-1} c_i^{(k)} A^i r, \quad c_i^{(k)} \in \mathbb{R}.$$

Как следствие, вектор $p_k = A^k r$ принадлежит подпространству, являющемуся линейной оболочкой векторов $r, Ar, \dots, A^k r$, или, что то же

самое, линейной оболочкой векторов p_0, p_1, \dots, p_k . По этой причине p_{k+1} выражается через предшествующие векторы как

$$p_{k+1} = Ap_k - \gamma_0^{(k)} p_0 - \dots - \gamma_k^{(k)} p_k$$

с какими-то коэффициентами $\gamma_0^{(k)}, \dots, \gamma_k^{(k)}$.

Домножая скалярно полученное соотношение на векторы p_0, p_1, \dots, p_k и привлекая условие ортогональности вектора p_{k+1} всем p_0, p_1, \dots, p_k , получим

$$\gamma_j^{(k)} = \frac{\langle Ap_k, p_j \rangle}{\langle p_j, p_j \rangle}, \quad j = 0, 1, \dots, k.$$

Но при $j = 0, 1, \dots, k-2$ справедливо $\langle Ap_k, p_j \rangle = 0$, так как $\langle Ap_k, p_j \rangle = \langle p_k, Ap_j \rangle$, а вектор Ap_j есть линейная комбинация векторов p_0, p_1, \dots, p_{j+1} , каждый из которых ортогонален к p_k при $j+1 < k$, т. е. $j \leq k-2$.

Итак, из коэффициентов $\gamma_j^{(k)}$ ненулевыми остаются лишь два коэффициента

$$\alpha_k = \gamma_k^{(k)} = \frac{\langle Ap_k, p_k \rangle}{\langle p_k, p_k \rangle},$$

$$\beta_k = \gamma_{k-1}^{(k)} = \frac{\langle Ap_k, p_{k-1} \rangle}{\langle p_{k-1}, p_{k-1} \rangle}.$$

Далее,

$$\langle Ap_k, p_{k-1} \rangle = \langle p_k, Ap_{k-1} \rangle = \langle p_k, p_k + \alpha_{k-1} p_{k-1} + \beta_{k-1} p_{k-2} \rangle = \langle p_k, p_k \rangle,$$

и поэтому

$$\beta_k = \frac{\langle p_k, p_k \rangle}{\langle p_{k-1}, p_{k-1} \rangle}.$$

Это завершает доказательство теоремы. ■

3.8 Метод прогонки

До сих пор не делалось никаких дополнительных предположений о структуре нулевых и ненулевых элементов в матрице системы. Но для большого числа систем линейных уравнений, встречающихся в практике математического моделирования ненулевые элементы заполняют

матрицу не полностью, образуя в ней те или иные правильные структуры — ленты, блоки, их комбинации и т. п. Естественно попытаться использовать это обстоятельство при конструировании более эффективных численных методов для решения СЛАУ с такими матрицами.

Метод прогонки, предложенный в 1952–53 годах И.М. Гельфандом и О.В. Локуциевским, предназначен для решения линейных систем уравнений с трёхдиагональными матрицами.¹⁸ Это важный в приложениях случай СЛАУ, возникающий, к примеру, при решении многих краевых задач для дифференциальных уравнений. По определению, трёхдиагональными называются матрицы, все ненулевые элементы которых сосредоточены на трёх диагоналях — главной и соседних с ней сверху и снизу. Иными словами, для трёхдиагональной матрицы $A = (a_{ij})$ равенство $a_{ij} \neq 0$ имеет место лишь при $i = j$ и $i = j \pm 1$.

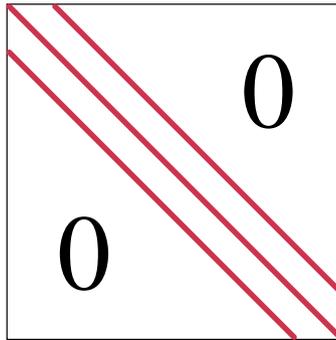


Рис. 3.18. Портрет трёхдиагональной матрицы.

Обычно трёхдиагональную систему n линейных уравнений с n неизвестными x_1, x_2, \dots, x_n записывают в следующем специальном каноническом виде, даже без обращения к матрично-векторной форме:

$$a_i x_{i-1} + b_i x_i + c_i x_{i+1} = d_i, \quad 1 \leq i \leq n, \quad (3.84)$$

где для единообразия рассматривают дополнительные фиктивные переменные x_0 и x_{n+1} и полагают $a_1 = c_n = 0$. Подобный вид и обозначения оправдываются тем, что соответствующие СЛАУ получают действительно «локально», как дискретизация дифференциальных

¹⁸Для краткости можно называть их просто трёхдиагональными линейными системами.

уравнений, связывающих значения искомых величин также локально, в окрестности какой-либо рассматриваемой точки.

Пример 3.8.1 В §2.8 мы могли видеть, что на равномерной сетке

$$u''(x_i) \approx \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2},$$

и правая часть этой формулы помимо самого узла x_i , в котором берётся производная, вовлекает ещё только соседние узлы x_{i-1} и x_{i+1} . Поэтому решение конечно-разностными методами краевых задач для различных дифференциальных уравнений второго порядка приводит к линейным системам уравнений с матрицами, у которых помимо главной диагонали заполнены только две соседние с ней, т. е. к системам с трёхдиагональными матрицами. ■

Соотношения вида (3.84)

$$a_i x_{i-1} + b_i x_i + c_i x_{i+1} = d_i, \quad i = 1, 2, \dots,$$

называют также *трёхточечными разностными уравнениями* или *разностными уравнениями второго порядка*.

Пусть для СЛАУ с трёхдиагональной матрицей выполняется прямой ход метода Гаусса без выбора ведущего элемента (т. е. без перестановок строк и столбцов матрицы). Если он успешно прорабатывает до конца, то приводит к системе с двухдиагональной матрицей вида

$$\begin{pmatrix} \times & \times & & \mathbf{0} \\ & \times & \ddots & \\ & & \ddots & \times \\ \mathbf{0} & & & \times & \times \\ & & & & \times \end{pmatrix}, \quad (3.85)$$

в которой ненулевые элементы присутствуют лишь на главной диагонали и первой верхней побочной. Следовательно, формулы обратного хода метода Гаусса вместо (3.54) должны иметь следующий двучленный вид

$$x_i = \xi_{i+1} x_{i+1} + \eta_{i+1}, \quad i = n, n-1, \dots, 1, \quad (3.86)$$

где, как и в исходных уравнениях, в n -ом соотношении присутствует вспомогательная фиктивная неизвестная x_{n+1} . Оказывается, что величины ξ_i и η_i в соотношениях (3.86) можно несложным образом выразить через элементы исходной системы уравнений.

Уменьшим в (3.86) все индексы на единицу —

$$x_{i-1} = \xi_i x_i + \eta_i$$

— и подставим полученное соотношение в i -ое уравнение системы, получим

$$a_i(\xi_i x_i + \eta_i) + b_i x_i + c_i x_{i+1} = d_i.$$

Отсюда

$$x_i = -\frac{c_i}{a_i \xi_i + b_i} x_{i+1} + \frac{d_i - a_i \eta_i}{a_i \xi_i + b_i}.$$

Сравнивая эту формулу с двучленными расчётными формулами (3.86), можем заключить, что

$$\xi_{i+1} = -\frac{c_i}{a_i \xi_i + b_i}, \quad (3.87)$$

$$\eta_{i+1} = \frac{d_i - a_i \eta_i}{a_i \xi_i + b_i}, \quad (3.88)$$

для $i = 1, 2, \dots, n$. Это формулы *прямого хода* прогонки, целью которого является вычисление величин ξ_i и η_i , называемых *прогоночными коэффициентами*. Вместе с формулами *обратного хода* (3.86) они определяют метод прогонки для решения систем линейных алгебраических уравнений с трёхдиагональной матрицей.

Для начала расчётов требуется знать величины ξ_1 и η_1 в прямом ходе и x_{n+1} — в обратном. Формально они неизвестны, но фактически полностью определяются условием $a_1 = c_n = 0$. Действительно, конкретные значения ξ_1 и η_1 не влияют на результаты решения, потому что в формулах (3.87)–(3.88) прямого хода прогонки они встречаются с множителем $a_1 = 0$. Кроме того, из формул прямого хода следует, что

$$\xi_{n+1} = -\frac{c_n}{a_n \xi_n + b_n} = -\frac{0}{a_n \xi_n + b_n} = 0,$$

а это коэффициент при x_{n+1} в обратном ходе прогонки. Поэтому и x_{n+1} может быть произвольным. Итак, для начала прогонки можно положить, к примеру,

$$\xi_1 = \eta_1 = x_{n+1} = 0. \quad (3.89)$$

Дадим теперь достаточные условия выполнимости метода прогонки, т. е. того, что знаменатели в расчётных формулах прямого хода не обращаются в нуль. Эти условия, фактически, будут также обосновывать возможность приведения трёхдиагональной матрицы исходной СЛАУ к двухдиагональному виду (3.85) преобразованиями прямого хода метода Гаусса без перестановки строк или столбцов, так как эти преобразования являются ничем иным, как прямым ходом метода прогонки.

Предложение 3.8.1 *Если в системе линейных алгебраических уравнений с трёхдиагональной матрицей (3.84) имеет место диагональное преобладание, т. е.*

$$|b_i| > |a_i| + |c_i|, \quad i = 1, 2, \dots, n,$$

то метод прогонки с выбором начальных значений согласно (3.89) является реализуемым.

По поводу формулировки Предложения 3.8.1 можно заметить, что условие диагонального преобладания в матрице влечёт её строгую регулярность, как мы видели в §3.6д. Поэтому в силу Теоремы 3.6.2 существует LU-разложение такой матрицы, и оно может быть получено с помощью прямого хода метода Гаусса без перестановки строк и столбцов. Но это и означает реализуемость метода прогонки. Ниже, тем не менее, даётся другое доказательство этого факта, которое позволяет помимо установления реализуемости дать ещё числовые оценки «запаса устойчивости» прогонки, т. е. того, насколько сильно знаменатели выражений (3.87)–(3.88) для прогоночных коэффициентов отличны от нуля в зависимости от элементов матрицы СЛАУ.

Доказательство. Покажем по индукции, что в рассматриваемой реализации прогонки для всех индексов i справедливо неравенство $|\xi_i| < 1$.

Прежде всего, $\xi_1 = 0$ и потому база индукции выполнена: $|\xi_1| < 1$. Далее, предположим, что для некоторого индекса i уже установлена оценка $|\xi_i| < 1$. Если соответствующее $c_i = 0$, то из (3.87) следует $\xi_{i+1} = 0$, и индукционный переход доказан. Поэтому пусть $c_i \neq 0$.

Тогда справедлива следующая цепочка соотношений

$$\begin{aligned}
 |\xi_{i+1}| &= \left| -\frac{c_i}{a_i \xi_i + b_i} \right| = \frac{|c_i|}{|a_i \xi_i + b_i|} \\
 &\leq \frac{|c_i|}{||b_i| - |a_i| \cdot |\xi_i||} \quad \text{из оценки снизу для модуля суммы} \\
 &< \frac{|c_i|}{|a_i| + |c_i| - |a_i| \cdot |\xi_i|} \quad \text{в силу диагонального преобладания} \\
 &= \frac{|c_i|}{|a_i|(1 - |\xi_i|) + |c_i|} \leq \frac{|c_i|}{|c_i|} = 1,
 \end{aligned}$$

где при переходе ко второй строке мы воспользовались известным неравенством для модуля суммы двух чисел:

$$|x + y| \geq ||x| - |y||. \quad (3.90)$$

Итак, неравенства $|\xi_i| < 1$ доказаны для всех прогоночных коэффициентов ξ_i , $i = 1, 2, \dots, n + 1$.

Как следствие, для знаменателей прогоночных коэффициентов ξ_i и η_i в формулах (3.87)–(3.88) имеем

$$\begin{aligned}
 |a_i \xi_i + b_i| &\geq ||b_i| - |a_i \xi_i|| \quad \text{по неравенству (3.90)} \\
 &= |b_i| - |a_i| |\xi_i| \quad \text{в силу диагонального преобладания} \\
 &> |a_i| + |c_i| - |a_i| \cdot |\xi_i| \quad \text{из-за диагонального преобладания} \\
 &= |a_i|(1 - |\xi_i|) + |c_i| \\
 &\geq |c_i| \geq 0 \quad \text{в силу оценки } |\xi_i| < 1,
 \end{aligned}$$

то есть строгое отделение от нуля. Это и требовалось доказать. ■

Отметим, что существуют и другие условия реализуемости метода прогонки. Например, некоторые из них требуют от матрицы «более мягкое» нестрогое диагональное преобладание (3.14), но зато более жёсткие, чем в Предложении 3.8.1, условия на коэффициенты системы. Весьма популярна, в частности, такая формулировка [17]:

Предложение 3.8.2 *Если в трёхдиагональной системе линейных алгебраических уравнений (3.84) побочные диагонали не содержат нулей,*

т. е. $a_i \neq 0$, $i = 2, 3, \dots, n$, и $c_i \neq 0$, $i = 1, 2, \dots, n - 1$, имеет место нестрогое диагональное преобладание

$$|b_i| \geq |a_i| + |c_i|, \quad i = 1, 2, \dots, n,$$

но хотя бы для одного индекса i это неравенство является строгим, то метод прогонки реализуем.

Нетрудно убедиться, что реализация прогонки требует линейного в зависимости от размера системы количества арифметических операций (примерно $8n$), т. е. весьма экономична.

На сегодняшний день разработано немало различных модификаций метода прогонки, которые хорошо приспособлены для решения тех или иных специальных систем уравнений, как трёхдиагональных, так и более общих, имеющих ленточные или даже блочно-ленточные матрицы [17]. В частности, существует метод матричной прогонки [27].

3.9 Стационарные итерационные методы для решения линейных систем

3.9а Краткая теория

Итерационные методы решения уравнений и систем уравнений — это методы, порождающие последовательность приближений $\{x^{(k)}\}_{k=0}^{\infty}$ к искомому решению x^* , которое получается как предел

$$x^* = \lim_{k \rightarrow \infty} x^{(k)}.$$

Допуская некоторую вольность речи, обычно говорят, что «итерационный метод сходится», если к пределу сходится конструируемая им последовательность приближений $\{x^{(k)}\}$.

Естественно, что на практике переход к пределу по $k \rightarrow \infty$ невозможен в силу конечности объема вычислений, который мы можем произвести. Поэтому при реализации итерационных методов вместо x^* обычно довольствуются нахождением какого-то достаточно хорошего приближения $x^{(k)}$ к x^* . Здесь важно правильно выбрать условие остановки итераций, при котором мы прекращаем порождать очередные приближения и выдаём $x^{(k)}$ в качестве решения. Подробнее мы рассмотрим этот вопрос в §3.14.

Общая схема итерационных методов выглядит следующим образом: выбирают одно или несколько *начальных приближений* $x^{(0)}, x^{(1)}, \dots, x^{(m)}$, а затем по их известным значениям последовательно вычисляются

$$x^{(k+1)} \leftarrow T_k(x^{(0)}, x^{(1)}, \dots, x^{(k)}), \quad k = m, m+1, m+2, \dots, \quad (3.91)$$

где T_k — отображение, называемое *оператором перехода* или *оператором шага* (k -го). Конечно, в реальных итерационных процессах каждое следующее приближение, как правило, зависит не от всех предшествующих приближений, а лишь от какого-то их фиксированного конечного числа. Более точно, итерационный процесс (3.91) называют *p -шаговым*, если его последующее приближение $x^{(k+1)}$ является функцией только от p предшествующих приближений, т. е. от $x^{(k)}, x^{(k-1)}, \dots, x^{(k-p+1)}$. В частности, наиболее простыми в этом отношении являются *одношаговые* итерационные методы

$$x^{(k+1)} \leftarrow T_k(x^{(k)}), \quad k = 0, 1, 2, \dots,$$

в которых $x^{(k+1)}$ зависит лишь от значения одной предшествующей итерации $x^{(k)}$. Для начала работы одношаговых итерационных процессов нужно знать одно начальное приближение $x^{(0)}$.

Итерационный процесс называется *стационарным*, если оператор перехода T_k не зависит от номера шага k , т. е. $T_k = T$, и *нестационарным* в противном случае. *Линейным p -шаговым итерационным процессом* будут называться итерации, в которых оператор перехода имеет вид

$$\begin{aligned} T_k(x^{(k)}, x^{(k-1)}, \dots, x^{(k-p+1)}) \\ = C^{(k,k)}x^{(k)} + C^{(k,k-1)}x^{(k-1)} + \dots + C^{(k,k-p+1)}x^{(k-p+1)} + d^{(k)} \end{aligned}$$

с какими-то коэффициентами $C^{(k,k)}, C^{(k,k-1)}, \dots, C^{(k,k-p+1)}$ и свободным членом $d^{(k)}$. В случае векторной неизвестной переменной x все $C^{(k,l)}$ являются матрицами подходящих размеров, а $d^{(k)}$ — вектор той же размерности, что и x . Матрицы $C^{(k,l)}$ часто называют *матрицами перехода* рассматриваемого итерационного процесса.

Итерационные методы были представлены выше в абстрактной манере, как некоторые конструктивные процессы, которые порождают последовательности, сходящиеся к искомому решению. В действительности, мотивации возникновения и развития итерационных методов

являлись существенно более ясными и практичными. Итерационные методы решения уравнений и систем уравнений возникли как уточняющие процедуры, которые позволяли за небольшое (удовлетворяющее практику) количество шагов получить приемлемое по точности приближённое решение задачи. Многие из классических итерационных методов явно несут отпечаток этих взглядов и ценностей.

Ясно, что для коррекции приближённого решения необходимо знать, насколько и как именно оно нарушает точное равенство обеих частей уравнения. На этом пути возникает важное понятие *невязки* приближённого решения \tilde{x} , которая определяется как разность левой и правой частей уравнения (системы уравнений) после подстановки в него \tilde{x} . Исследование этой величины, отдельных её компонент (в случае системы уравнений) и решение вопроса о том, как можно на основе этой информации корректировать приближение к решению, составляет важнейшую часть работы по конструированию итерационных методов.

Мы подробно рассматриваем различные итерационные методы для решения нелинейных уравнений и систем уравнений в Главе 4, а здесь основное внимание будет уделено итерационному решению систем линейных алгебраических уравнений и проблемы собственных значений.

Причины, по которым для решения систем линейных уравнений итерационные методы могут оказаться более предпочтительными, чем прямые, заключаются в следующем. Большинство итерационных методов являются *самоисправляющимися*, т. е. такими, в которых погрешность, допущенная в вычислениях, при сходимости исправляется в ходе итерирования и не отражается на окончательном результате. Это следует из конструкции оператора перехода, в котором обычно по самому его построению присутствует информация о решаемой системе уравнений (что мы увидим далее на примерах). При выполнении алгоритма эта информация на каждом шаге вносится в итерационный процесс и оказывает влияние на его ход. Напротив, прямые методы решения СЛАУ этим свойством не обладают, так как, оттолкнувшись от исходной системы, мы далее уже не возвращаемся к ней, а оперируем с её следствиями, которые никакой обратной связи от исходной системы не получают.¹⁹

Нередко итерационные процессы сравнительно несложно программируются, так как представляют собой повторяющиеся единообразные

¹⁹Для исправления этого положения прямые методы решения СЛАУ в ответственных ситуациях часто дополняют процедурами итерационного уточнения. См., к примеру, пункт 67 главы 4 в [42].

процедуры, применяемые к последовательным приближениям к решению. При решении СЛАУ с разреженными матрицами в итерационных процессах нередко можно легче, чем в прямых методах, учитывать структуру нулевых и ненулевых элементов матрицы и основывать на этом упрощённые формулы матрично-векторного умножения, которые существенно уменьшают общую трудоёмкость алгоритма.

Иногда системы линейных алгебраических уравнений задаются в операторном виде, рассмотренном нами в начале §3.6 (стр. 275) т. е. так, что их матрица и правая часть не выписываются явно. Вместо этого задаётся действие такой матрицы (линейного оператора) на любой вектор, и это позволяет строить и использовать итерационные методы. С другой стороны, преобразования матриц таких систем, которые являются основой прямых методов решения систем линейных уравнений, очень сложны или порой просто невозможны.

Наконец, быстро сходящиеся итерационные методы могут обеспечивать выигрыш по времени даже для СЛАУ общего вида, если требуют небольшое число итераций.

То обстоятельство, что искомое решение получается как топологический предел последовательности, порождаемой методом, является характерной чертой именно итерационных методов решения уравнений. Существуют и другие конструкции, по которым решение задачи строится из последовательности, порождаемой методом. Интересный пример дают методы Монте-Карло, в которых осуществляется усреднение последовательности приближений.

3.9б Сходимость стационарных одношаговых итерационных методов

Системы линейных уравнений вида

$$x = Cx + d,$$

в котором вектор неизвестных переменных выделен в одной из частей, мы будем называть *системами в рекуррентном виде*.

Теорема 3.9.1 Пусть система уравнений $x = Cx + d$ имеет единственное решение. Стационарный одношаговый итерационный процесс

$$x^{(k+1)} \leftarrow Cx^{(k)} + d, \quad k = 0, 1, 2, \dots, \quad (3.92)$$

сходится при любом начальном приближении $x^{(0)}$ тогда и только тогда, когда спектральный радиус матрицы C меньше единицы, т. е. $\rho(C) < 1$.

Оговорка о единственности решения существенна. Если взять, к примеру, $C = I$ и $d = 0$, то рассматриваемая система обратится в тождество $x = x$, имеющее решением любой вектор. Соответствующий итерационный процесс $x^{(k+1)} \leftarrow x^{(k)}$, $k = 0, 1, 2, \dots$, будет сходиться из любого начального приближения, хотя спектральный радиус матрицы перехода C равен единице.

Доказательство Теоремы 3.9.1 будет разбито на две части, результат каждой из которых представляет самостоятельный интерес.

Предложение 3.9.1 Если $\|C\| < 1$ в какой-нибудь матричной норме, то стационарный одношаговый итерационный процесс

$$x^{(k+1)} \leftarrow Cx^{(k)} + d, \quad k = 0, 1, 2, \dots,$$

сходится при любом начальном приближении $x^{(0)}$.

Доказательство. В формулировке предложения ничего не говорится о пределе, к которому сходится последовательность приближений $\{x^{(k)}\}$, порождаемых итерационным процессом. Но мы можем указать его в явном виде и строить доказательство с учётом этого знания.

Если $\|C\| < 1$ для какой-нибудь матричной нормы, то в силу результата о матричном ряде Неймана (Предложение 3.3.11, стр. 254) матрица $(I - C)$ неособенна и имеет обратную. Следовательно, система уравнений $(I - C)x = d$, как и равносильная ей $x = Cx + d$, имеют единственное решение, которое мы обозначим x^* . Покажем, что в условиях предложения это и есть предел последовательных приближений $x^{(k)}$.

В самом деле, если

$$x^* = Cx^* + d,$$

то, вычитая это равенство из соотношений $x^{(k)} = Cx^{(k-1)} + d$, $k = 1, 2, \dots$, получим

$$x^{(k)} - x^* = C(x^{(k-1)} - x^*).$$

Вспомним, что всякая матричная норма согласована с некоторой векторной нормой (Предложение 3.3.4), и именно эту норму мы применим к обеим частям последнего равенства:

$$\|x^{(k)} - x^*\| \leq \|C\| \|x^{(k-1)} - x^*\|.$$

Повторное применение этой оценки погрешности для $x^{(k-1)}, x^{(k-2)}, \dots$, и т. д. вплоть до $x^{(1)}$ приводит к цепочке неравенств

$$\begin{aligned} \|x^{(k)} - x^*\| &\leq \|C\| \cdot \|x^{(k-1)} - x^*\| \\ &\leq \|C\|^2 \cdot \|x^{(k-2)} - x^*\| \\ &\leq \dots \quad \dots \\ &\leq \|C\|^k \cdot \|x^{(0)} - x^*\|. \end{aligned} \quad (3.93)$$

Правая часть неравенства (3.93) сходится к нулю при $k \rightarrow \infty$ в силу условия $\|C\| < 1$, поэтому последовательность приближений $\{x^{(k)}\}_{k=0}^{\infty}$ действительно сходится к пределу x^* . ■

Побочным следствием доказательства Предложения 3.9.1 является прояснение роли нормы матрицы перехода $\|C\|$ как коэффициента подавления погрешности приближений к решению СЛАУ в согласованной векторной норме. Это следует из неравенств (3.93): чем меньше $\|C\|$, тем быстрее убывает эта погрешность на каждом отдельном шаге итерационного процесса.

Предложение 3.9.2 *Для любой квадратной матрицы A и любого $\epsilon > 0$ существует такая подчинённая матричная норма $\|\cdot\|_{\epsilon}$, что*

$$\rho(A) \leq \|A\|_{\epsilon} \leq \rho(A) + \epsilon.$$

Доказательство. Левое из выписанных неравенств было обосновано ранее в Предложении 3.3.9, и потому содержанием сформулированного результата является правое неравенство, дающее, фактически, оценку снизу для спектрального радиуса с помощью некоторой специальной матричной нормы.

С помощью преобразования подобия приведём матрицу A к жордановой канонической форме

$$S^{-1}AS = J,$$

где

$$J = \left(\begin{array}{ccc|ccc} \lambda_1 & 1 & & & & \\ & \lambda_1 & \ddots & & & \\ & & \ddots & & & \\ & & & 1 & & \\ & & & \lambda_1 & & \\ \hline & & & & \lambda_2 & 1 & \\ & 0 & & & \ddots & \ddots & \\ & & & & & \lambda_2 & \\ \hline & & & & & & \\ & 0 & & & 0 & & \\ & & & & & & \ddots \\ & & & & & & \ddots \end{array} \right),$$

а S — некоторая неособенная матрица, осуществляющая преобразование подобия. Положим

$$D_\epsilon := \text{diag} \{1, \epsilon, \epsilon^2, \dots, \epsilon^{n-1}\}$$

— диагональной $n \times n$ -матрице с числами $1, \epsilon, \epsilon^2, \dots, \epsilon^{n-1}$ по главной диагонали. Тогда нетрудно проверить, что

$$(SD_\epsilon)^{-1}A(SD_\epsilon) = D_\epsilon^{-1}(S^{-1}AS)D_\epsilon$$

$$= D_\epsilon^{-1}JD_\epsilon = \left(\begin{array}{ccc|ccc} \lambda_1 & \epsilon & & & & \\ & \lambda_1 & \ddots & & & \\ & & \ddots & & & \\ & & & \epsilon & & \\ & & & \lambda_1 & & \\ \hline & & & & \lambda_2 & \epsilon & \\ & & & & \ddots & \ddots & \\ & & & & & \lambda_2 & \\ \hline & & & & & & \\ & & & & & & \ddots \\ & & & & & & \ddots \end{array} \right),$$

— матрица в «модифицированной» жордановой форме, которая отличается от обычной жордановой формы присутствием ϵ вместо 1 на верхней побочной диагонали каждой жордановой клетки.

Действительно, умножение на диагональную матрицу слева — это умножение строк матрицы на соответствующие диагональные элементы, а умножение на диагональную матрицу справа равносильно умножению столбцов на элементы диагонали. Два таких умножения — на $D_\epsilon^{-1} = \text{diag}\{1, \epsilon^{-1}, \epsilon^{-2}, \dots, \epsilon^{1-n}\}$ слева и на $D_\epsilon = \text{diag}\{1, \epsilon, \epsilon^2, \dots, \epsilon^{n-1}\}$ справа — компенсируют друг друга на главной диагонали матрицы J . Но на верхней побочной диагонали, где ненулевые элементы имеют индексы $(i, i-1)$, от этих умножений остаётся множитель $\epsilon^{-i}\epsilon^{i-1} = \epsilon$, $i = 0, 1, \dots, n-1$.

Определим теперь векторную норму

$$\|x\|_\epsilon := \|(SD_\epsilon)^{-1}x\|_\infty.$$

Тогда для подчинённой ей матричной нормы справедлива следующая цепочка оценок

$$\begin{aligned} \|A\|_\epsilon &= \max_{x \neq 0} \frac{\|Ax\|_\epsilon}{\|x\|_\epsilon} = \max_{x \neq 0} \frac{\|(SD_\epsilon)^{-1}Ax\|_\infty}{\|(SD_\epsilon)^{-1}x\|_\infty} \\ &= \max_{y \neq 0} \frac{\|(SD_\epsilon)^{-1}A(SD_\epsilon)y\|_\infty}{\|y\|_\infty} \quad \text{после замены } y = (SD_\epsilon)^{-1}x \\ &= \max_{y \neq 0} \frac{\|(D_\epsilon^{-1}JD_\epsilon)y\|_\infty}{\|y\|_\infty} = \|D_\epsilon^{-1}JD_\epsilon\|_\infty \\ &= \text{максимум сумм модулей элементов в } D_\epsilon^{-1}JD_\epsilon \text{ по строкам} \\ &\leq \max_i |\lambda_i(A)| + \epsilon = \rho(A) + \epsilon, \end{aligned}$$

где $\lambda_i(A)$ — i -ое собственное значение матрицы A . Неравенство при переходе к последней строке возникает по существу, так как матрица может иметь большое по модулю собственное значение в жордановой клетке размера 1×1 , в которой нет элементов ϵ . ■

Доказательство Теоремы 3.9.1 о сходимости одношагового стационарного итерационного процесса.

Сначала покажем необходимость условия теоремы. Пусть порождаемая в итерационном процессе последовательность $\{x^{(k)}\}$ сходится. Её пределом при этом может быть только решение x^* системы $x = Cx + d$,

т. е. должно быть $\lim_{k \rightarrow \infty} x^{(k)} = x^*$, в чём можно убедиться, переходя в соотношении

$$x^{(k+1)} = Cx^{(k)} + d$$

к пределу по $k \rightarrow \infty$. Далее, вычитая почленно равенство для точного решения $x^* = Cx^* + d$ из расчётной формулы итерационного процесса $x^{(k)} = Cx^{(k-1)} + d$, получим

$$x^{(k)} - x^* = C(x^{(k-1)} - x^*), \quad k = 1, 2, \dots,$$

откуда

$$\begin{aligned} x^{(k)} - x^* &= C(x^{(k-1)} - x^*) \\ &= C^2(x^{(k-2)} - x^*) \\ &= \dots \dots \\ &= C^k(x^{(0)} - x^*). \end{aligned}$$

Так как левая часть этих равенств при $k \rightarrow \infty$ сходится к нулю, то должна сходиться к нулю и правая, причём для любого вектора $x^{(0)}$. В силу единственности и, как следовательно, фиксированности решения x^* вектор $(x^{(0)} - x^*)$ также может быть произвольным, и тогда сходимость погрешности к нулю возможна лишь при $C^k \rightarrow 0$. На основании Предложения 3.3.10 (стр. 252) заключаем, что спектральный радиус C должен быть строго меньше 1.

Достаточность. Если $\rho(C) < 1$, то взяв положительное ϵ удовлетворяющим оценке $\epsilon < 1 - \rho(C)$, мы можем согласно Предложению 3.9.2 выбрать матричную норму $\|\cdot\|_\epsilon$ так, чтобы выполнялось неравенство $\|C\|_\epsilon < 1$. Далее в этих условиях применимо Предложение 3.9.1, которое утверждает сходимость итерационного процесса (3.92)

$$x^{(k+1)} \leftarrow Cx^{(k)} + d, \quad k = 0, 1, 2, \dots$$

Это завершает доказательство Теоремы 3.9.1. ■

Доказанные результаты — теорема и два предложения — проясняют роль спектрального радиуса среди различных характеристик матрицы. Мы могли видеть в §3.3ж, что спектральный радиус не является матричной нормой, но, как выясняется, его с любой степенью точности

можно приблизить некоторой подчинённой матричной нормой. Кроме того, понятие спектрального радиуса оказывается чрезвычайно полезным при исследовании итерационных процессов и вообще степеней матрицы.

Следствие из Предложения 3.9.2. Степени матрицы A^k сходятся к нулевой матрице при $k \rightarrow \infty$ тогда и только тогда, когда $\rho(A) < 1$.

В самом деле, ранее мы установили (Предложение 3.3.10), что из сходимости степеней матрицы A^k при $k \rightarrow \infty$ к нулевой матрице вытекает $\rho(A) < 1$. Теперь результат Предложения 3.9.2 позволяет сказать, что это условие на спектральный радиус является и достаточным: если $\rho(A) < 1$, то мы можем подобрать матричную норму так, чтобы $\|A\| < 1$, и тогда $\|A^n\| \leq \|A\|^n \rightarrow 0$ при $n \rightarrow \infty$.

С учётом Предложения 3.9.2 более точно переформулируются условия сходимости матричного ряда Неймана (Предложение 3.3.11): он сходится для матрицы A тогда и только тогда, когда $\rho(A) < 1$, а условие $\|A\| < 1$ является всего лишь достаточным.

Заметим, что для несимметричных матриц нормы, близкие к спектральному радиусу, могут оказаться очень экзотичными и даже неестественными. Это видно из доказательства Теоремы 3.9.1. Как правило, исследовать сходимость итерационных процессов лучше всё-таки в обычных нормах, часто имеющих практический смысл.

Интересен вопрос о выборе начального приближения для итерационных методов решения СЛАУ. Иногда его решают из каких-то содержательных соображений, когда в силу физических и прочих содержательных причин бывает известно некоторое хорошее приближение к решению, а итерационный метод предназначен для его уточнения. При отсутствии таких условий начальное приближение нужно выбирать на основе других идей.

Например, если в рекуррентном виде $x = Cx + d$, исходя из которого строятся сходящиеся итерации, матрица C имеет «малую» норму (относительно неё мы вправе предполагать, что $\|C\| < 1$), то тогда членом Cx можно пренебречь. Как следствие, точное решение не сильно отличается от вектора свободных членов d , и поэтому можно взять $x^{(0)} = d$. Этот вектор привлекателен также тем, что получается как первая итерация при нулевом начальном приближении. Беря $x^{(0)} = d$, мы экономим на этой итерации.

3.9в Подготовка линейной системы к итерационному процессу

В этом параграфе мы исследуем различные способы приведения системы линейных алгебраических уравнений

$$Ax = b \quad (3.94)$$

к равносильной системе в рекуррентном виде

$$x = Cx + d, \quad (3.95)$$

отталкиваясь от которого можно организовывать одношаговый итерационный процесс для решения (3.94). Фактически, это вопрос о том, как связан предел стационарного одношагового итерационного процесса (3.92) с интересующим нас решением системы линейных алгебраических уравнений $Ax = b$. При этом практический интерес представляет, естественно, не всякое приведение системы (3.94) к виду (3.95), но лишь такое, которое удовлетворяет условию сходимости стационарного одношагового итерационного процесса, выведенному в предшествующем разделе, т. е. $\rho(C) < 1$.

Существует большое количество различных способов приведения исходной ИСЛАУ к виду, допускающему применение итераций, большое разнообразие способов организации этих итерационных процессов и т. п. Не претендуя на всеохватную теорию, мы рассмотрим ниже лишь несколько общих приёмов подготовки и организации итерационных процессов.

Простейший способ состоит в том, чтобы добавить к обеим частям исходной системы по вектору неизвестной переменной x , т. е.

$$x + Ax = x + b, \quad (3.96)$$

а затем член Ax перенести в правую часть:

$$x = (I - A)x + b.$$

Иногда этот приём работает, но весьма часто он непригоден, так как спектральный радиус матрицы $C = I - A$ оказывается не меньшим единицы.

В самом деле, если λ — собственное значение для A , то для матрицы $(I - A)$ собственным значением будет $1 - \lambda$, и тогда $1 - \lambda > 1$ при

вещественных отрицательных λ . С другой стороны, если у матрицы A есть собственные значения, большие по модулю, чем 2, т. е. если $|\lambda| > 2$, то $|1 - \lambda| = |\lambda - 1| \geq ||\lambda| - 1| > 1$ и сходимости стационарных итераций мы также не получим.

Из предшествующих рассуждений можно ясно видеть, что необходим активный способ управления свойствами матрицы C в получающейся системе рекуррентного вида $x = Cx + d$. Одним из важнейших инструментов такого управления служит *предобуславливание* исходной системы.

Определение 3.9.1 *Предобуславливанием системы линейных алгебраических уравнений $Ax = b$ называется умножение слева обеих её частей на некоторую матрицу L . Сама эта матрица L называется предобуславливающей матрицей или, коротко, предобуславливателем.*

Цель предобуславливания — изменение (вообще говоря, улучшение) свойств матрицы A исходной системы $Ax = b$, вместо которой мы получаем систему

$$(LA)x = Lb.$$

Продуманный выбор предобуславливателя может, к примеру, изменить выгодным нам образом расположение спектра матрицы A , так необходимое для организации сходящихся итерационных процессов.

Естественно выполнить предобуславливание до перехода к системе (3.96), т. е. до прибавления вектора неизвестных x к обеим частям исходной СЛАУ. Поскольку тогда вместо $Ax = b$ будем иметь $(LA)x = Lb$, то далее получаем

$$x = (I - LA)x + Lb.$$

Теперь в этом рекуррентном виде с помощью подходящего выбора L можно добиваться требуемых свойств матрицы $(I - LA)$.

Каким образом следует выбирать предобуславливатели? Совершенно общего рецепта на этот счёт не существует, и теория разбивается здесь на набор рекомендаций для ряда более или менее конкретных важных случаев.

Например, если в качестве предобуславливающей матрицы взять $L = A^{-1}$ или хотя бы приближённо равную обратной к A , то вместо системы $Ax = b$ получим $(A^{-1}A)x = A^{-1}b$, т. е. систему уравнений

$$Ix = A^{-1}b$$

или близкую к ней, матрица которой обладает всеми возможными достоинствами (хорошим диагональным преобладанием, малой обусловленностью и т. п.). Ясно, что нахождение подобного предобуславливателя не менее трудно, чем решение исходной системы, но сама идея примера весьма плодотворна. На практике в качестве предобуславливателей часто берут несложно вычисляемые обратные матрицы к той или иной «существенной» части матрицы A . Например, к главной диагонали матрицы или же к трём диагоналям — главной и двум побочным.

Другой способ приведения СЛАУ к рекуррентному виду основан на *расщеплении* матрицы системы.

Определение 3.9.2 *Расщеплением квадратной матрицы A называется её представление в виде $A = G + (-H) = G - H$, где G — неособенная матрица.*

Если известно некоторое расщепление матрицы A , $A = G - H$, то вместо исходной системы $Ax = b$ мы можем рассмотреть

$$(G - H)x = b,$$

которая равносильна

$$Gx = Hx + b,$$

так что

$$x = G^{-1}Hx + G^{-1}b.$$

На основе полученного рекуррентного вида можно организовать итерации

$$x^{(k+1)} \leftarrow G^{-1}Hx^{(k)} + G^{-1}b, \quad (3.97)$$

задавшись каким-то начальным приближением $x^{(0)}$. Таким образом, всякое расщепление матрицы СЛАУ помогает конструированию итерационных процессов.

Но практическое значение имеют не все расщепления, а лишь те, в которых матрица G обращается «относительно просто», чтобы организация итерационного процесса не сделалась более сложной задачей, чем решение исходной СЛАУ. Другое требование к матрицам, образующим расщепление, состоит в том, чтобы норма обратной для G , т. е. $\|G^{-1}\|$, была «достаточно малой», поскольку $\|G^{-1}H\| \leq \|G^{-1}\| \|H\|$. Если G^{-1} имеет большую норму, то может оказаться $\rho(G^{-1}H) > 1$, и для итерационного процесса (3.97) не будут выполнены условия сходимости.

Очень популярный способ расщепления матрицы A состоит в том, чтобы сделать элементы в $G = (g_{ij})$ и $H = (h_{ij})$ взаимнодополнительными, т. е. такими, что $g_{ij}h_{ij} = 0$ для любых индексов i и j . Тогда ненулевые элементы матриц G и $(-H)$ совпадают с ненулевыми элементами A .

В качестве примеров несложно обрацаемых матриц можно указать²⁰

- 1) диагональные матрицы,
- 2) треугольные матрицы,
- 3) трёхдиагональные матрицы,
- 4)

Ниже в §3.9д и §3.9е мы подробно рассмотрим итерационные процессы, соответствующие первым двум пунктам этого списка.

3.9г Скалярный преобуславливатель и его оптимизация

Напомним, что *скалярными матрицами* (из-за своего родства скалярам) называются матрицы, кратные единичным, т. е. имеющие вид τI , где $\tau \in \mathbb{R}$ или \mathbb{C} . Сейчас мы подробно исследуем описанную в предыдущем разделе возможность управления итерационным процессом на простейшем примере преобуславливания с помощью скалярной матрицы, когда $A = \tau I$, $\tau \in \mathbb{R}$ и $\tau \neq 0$.

Итак, рассматриваем итерационный процесс

$$x^{(k+1)} \leftarrow (I - \tau A)x^{(k)} + \tau b, \quad (3.98)$$

$\tau = \text{const}$, который часто называют *методом простой итерации*. Если λ_i , $i = 1, 2, \dots, n$, — собственные числа матрицы A (вообще говоря, они комплексны), то собственные числа матрицы $(I - \tau A)$ равны $(1 - \tau\lambda_i)$. Ясно, что в случае, когда среди λ_i имеются числа с разным знаком вещественной части $\text{Re } \lambda_i$, при любом фиксированном вещественном τ выражение

$$\text{Re}(1 - \tau\lambda_i) = 1 - \tau \text{Re } \lambda_i$$

будет иметь как меньшие 1 значения для каких-то λ_i , так и большие чем 1 значения для некоторых других λ_i . Как следствие, добиться ло-

²⁰Обратная матрица очень просто находится также для ортогональных матриц, но они не очень хороши для расщепления, так как норма обратной для них не мала.

кализации всех значений $(1 - \tau\lambda_i)$ в единичном круге комплексной плоскости с центром в нуле, т. е. соблюдения условия $\rho(I - \tau A) < 1$, никаким выбором τ будет невозможно.

Далее рассмотрим практически важный частный случай, когда A — симметричная положительно определённая матрица, так что все λ_i , $i = 1, 2, \dots, n$, вещественны и положительны. Обычно они не бывают известными, но нередко более или менее точно известен интервал их расположения на вещественной оси \mathbb{R} . Будем предполагать, что $\lambda_i \in [\mu, M]$, $i = 1, 2, \dots, n$.

Матрица $(I - \tau A)$ тогда также симметрична, и потому её спектральный радиус совпадает с 2-нормой. Чтобы обеспечить сходимость итерационного процесса и добиться её наибольшей скорости, нам нужно, согласно Теореме 3.9.1 и оценкам убывания погрешности (3.93), найти значение τ , которое доставляет минимум величине

$$\|I - \tau A\|_2 = \max_{\lambda_i} |1 - \tau\lambda_i|,$$

где максимум в правой части берётся по дискретному множеству точек λ_i , $i = 1, 2, \dots, n$, спектра матрицы A . В условиях, когда о расположении λ_i ничего не известно кроме их принадлежности интервалу $[\mu, M]$, естественно заменить максимизацию по множеству λ_i , $i = 1, 2, \dots, n$, на максимизацию по всему объемлющему его интервалу $[\mu, M]$. Итак, мы будем искать оптимальное значение $\tau = \tau_{\text{опт}}$, при котором достигается

$$\min_{\tau} \left(\max_{\lambda \in [\mu, M]} |1 - \tau\lambda| \right).$$

Обозначив

$$g(\tau) = \max_{\mu \leq \lambda \leq M} |1 - \tau\lambda|,$$

обратимся для минимизации функции $g(\tau)$ к геометрической иллюстрации Рис. 3.19. Пользуясь ею, мы исследуем поведение $g(\tau)$ при изменении аргумента τ .

При $\tau \leq 0$ функция $(1 - \tau\lambda)$ не убывает по λ , и при положительных λ , очевидно, не меньше 1 по абсолютной величине. Тогда итерационный процесс (3.98) сходиться не будет. Следовательно, в нашем анализе имеет смысл ограничиться теми τ , для которых $(1 - \tau\lambda)$ убывает по λ . Это значения $\tau > 0$.

При $0 < \tau \leq M^{-1}$ функция $1 - \tau\lambda$ на интервале $\lambda \in [\mu, M]$ неотрицательна и монотонно убывает. Поэтому $g(\tau) = \max_{\lambda} |1 - \tau\lambda| = 1 - \tau\mu$ и достигается на левом конце интервала $[\mu, M]$.

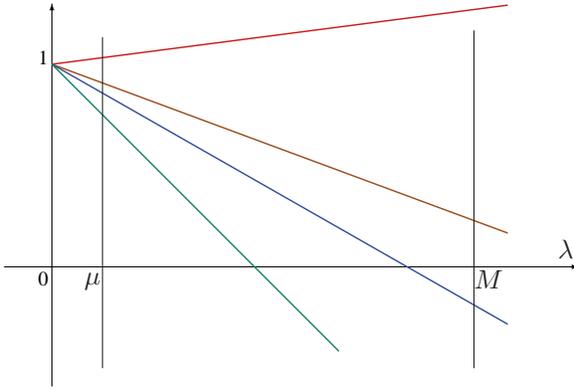


Рис. 3.19. Графики функций $1 - \tau\lambda$ для различных τ

При $\tau > M^{-1}$ величина $1 - \tau M$ отрицательна, так что график функции $1 - \tau\lambda$ на интервале $\lambda \in [\mu, M]$ пересекает ось абсцисс. При этом

$$g(\tau) = \max\{1 - \tau\mu, -(1 - \tau M)\},$$

причём на левом конце $(1 - \tau\mu)$ убывает с ростом τ , а на правом конце $-(1 - \tau M)$ растёт с ростом τ .

При некотором $\tau = \tau_{\text{опт}}$ наступает момент, когда эти значения на концах интервала $[\mu, M]$ сравниваются друг с другом:

$$1 - \tau\mu = -(1 - \tau M).$$

Он и является моментом достижения оптимума, поскольку дальнейшее увеличение τ приводит к росту $-(1 - \tau M)$ на правом конце интервала, а уменьшение τ ведёт к росту $1 - \tau\mu$ на левом конце. В любом из этих случаев $g(\tau)$ возрастает. Отсюда

$$\tau_{\text{опт}} = \frac{2}{M + \mu}, \tag{3.99}$$

а значение оптимума $g(\tau)$, равное коэффициенту подавления 2-нормы погрешности (как следствие из неравенств (3.93)), есть

$$\begin{aligned} \|I - \tau_{\text{опт}}A\|_2 &= \min_{\tau} \max_{\lambda \in [\mu, M]} |1 - \tau\lambda| = 1 - \tau_{\text{опт}}\mu \\ &= 1 - \frac{2}{M + \mu} \cdot \mu = \frac{M - \mu}{M + \mu}. \end{aligned} \tag{3.100}$$

Ясно, что эта величина меньше единицы, т. е. даже с помощью простейшего скалярного предобуславливателя мы добились сходимости итерационного процесса.

Полезно оценить значение (3.100), используя спектральное число обусловленности матрицы A . Так как $\mu \leq \lambda_{\min}(A)$ и $\lambda_{\max}(A) \leq M$, то

$$\text{cond}_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \leq \frac{M}{\mu}.$$

Поэтому, принимая во внимание тот факт, что функция

$$f(x) = \frac{x-1}{x+1} = 1 - \frac{2}{x+1}$$

возрастает при положительных x , можем заключить, что

$$\|I - \tau_{\text{онт}}A\|_2 = \frac{M/\mu - 1}{M/\mu + 1} \geq \frac{\text{cond}_2(A) - 1}{\text{cond}_2(A) + 1}.$$

Получается, что чем больше $\text{cond}_2(A)$, т. е. чем хуже обусловленность матрицы A исходной системы, тем медленнее сходимость нашего итерационного процесса. Мы увидим далее, что это характерно для поведения многих итерационных методов.

Наибольшую трудность на практике представляет нахождение μ , т. е. нижней границы спектра матрицы СЛАУ. Иногда мы даже можем ничего не знать о её конкретной величине кроме того, что $\mu \geq 0$. В этих условиях развитая нами теория применима лишь частично. Оптимальным значением параметра τ следует, очевидно, взять

$$\tau_{\text{онт}} = \frac{2}{M},$$

метод простой итерации (3.98) будет при этом сходиться, но никаких оценок его скорости сходимости дать нельзя.

3.9д Итерационный метод Якоби

Пусть в системе линейных алгебраических уравнений $Ax = b$ диагональные элементы матрицы $A = (a_{ij})$ отличны от нуля, т. е. $a_{ii} \neq 0$, $i = 1, 2, \dots, n$. Это условие не является обременительным, так как для неособенной матрицы A перестановкой строк (соответствующей перестановке уравнений системы) можно всегда сделать диагональные элементы ненулевыми.

В развёрнутом виде рассматриваемая система имеет вид

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

и, выражая i -ю компоненту вектора неизвестных из i -го уравнения, получим

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

Нетрудно понять, что эти соотношения дают представление исходной СЛАУ в рекуррентном виде $x = T(x)$, необходимом для организации одношаговых итераций $x^{(k+1)} \leftarrow T(x^{(k)})$, $k = 0, 1, 2, \dots$. Здесь

$$T(x) = (T_1(x), T_2(x), \dots, T_n(x))^T$$

и

$$T_i(x) = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij}x_j \right), \quad i = 1, 2, \dots, n.$$

Таблица 3.5. Итерационный метод Якоби для решения СЛАУ

```

k ← 0;
выбираем начальное приближение x(0);
DO WHILE ( метод не сошёлся )
  DO FOR i = 1 TO n
    xi(k+1) ←  $\frac{1}{a_{ii}} \left( b_i - \sum_{j \neq i} a_{ij}x_j^{(k)} \right)$ 
  END DO
  k ← k + 1;
END DO

```

Псевдокод соответствующего итерационного процесса представлен в Табл. 3.5, где вспомогательная переменная k — это счётчик числа итераций. Он был предложен ещё в середине XIX века К.Г. Якоби и часто (особенно в старых книгах по численным методам) называется «методом одновременных смещений». Под «смещениями» здесь имеются в виду коррекции компонент очередного приближения к решению, выполняемые на каждом шаге итерационного метода. Смещения-коррекции «одновременны» потому, что все компоненты следующего приближения $x^{(k+1)}$ насчитываются независимо друг от друга по единообразным формулам, основанным на использовании лишь предыдущего приближения $x^{(k)}$. В следующем параграфе будет рассмотрен итерационный процесс, устроенный несколько по-другому, в котором смещения-коррекции компонент очередного приближения к решению «не одновременны» в том смысле, что находятся последовательно одна за другой не только из предыдущего приближения, но ещё и друг из друга.

Пусть $A = \tilde{L} + D + \tilde{U}$, где

$$\tilde{L} = \begin{pmatrix} 0 & & & & \\ a_{21} & 0 & & & \\ a_{31} & a_{32} & \ddots & & \\ \vdots & \vdots & \ddots & & 0 \\ a_{n1} & a_{n2} & \cdots & a_{n,n-1} & 0 \end{pmatrix} \quad \begin{array}{l} \text{— строго нижняя} \\ \text{треугольная матрица.} \end{array}$$

$$D = \text{diag} \{a_{11}, a_{22}, \dots, a_{nn}\} \quad \text{— диагональ матрицы } A,$$

$$\tilde{U} = \begin{pmatrix} 0 & a_{12} & \cdots & a_{1,n-1} & a_{1n} \\ & 0 & \ddots & a_{2,n-1} & a_{2n} \\ & & \ddots & \vdots & \vdots \\ & & & 0 & a_{n-1,n} \\ \mathbf{0} & & & & 0 \end{pmatrix} \quad \begin{array}{l} \text{— строго верхняя} \\ \text{треугольная матрица.} \end{array}$$

Тогда итерационный метод Якоби может быть представлен как метод, основанный на таком расщеплении матрицы системы $A = G - H$ (см. §3.9в), что

$$G = D, \quad H = -(\tilde{L} + \tilde{U}).$$

Соответственно, в матричном виде метод Якоби записывается как

$$x^{(k+1)} \leftarrow -D^{-1}(\tilde{L} + \tilde{U})x^{(k)} + D^{-1}b, \quad k = 0, 1, 2, \dots$$

Теперь нетрудно дать условия его сходимости, основываясь на общем результате о сходимости стационарных одношаговых итераций (Теорема 3.9.1). Именно, метод Якоби сходится из любого начального приближения тогда и только тогда, когда

$$\rho(D^{-1}(\tilde{L} + \tilde{U})) < 1.$$

Матрица $D^{-1}(\tilde{L} + \tilde{U})$ просто выписывается по исходной системе и имеет вид

$$\begin{pmatrix} 0 & a_{12}/a_{11} & \dots & a_{1n}/a_{11} \\ a_{21}/a_{22} & 0 & \dots & a_{2n}/a_{22} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1}/a_{nn} & a_{n2}/a_{nn} & \dots & 0 \end{pmatrix}. \quad (3.101)$$

Но нахождение её спектрального радиуса является задачей, сравнимой по сложности с выполнением самого итерационного процесса, и потому применять его для исследования сходимости метода Якоби непрактично. Для быстрой и грубой оценки спектрального радиуса можно воспользоваться какой-нибудь матричной нормой и результатом Предложения 3.3.9.

Полезен также следующий достаточный признак сходимости:

Предложение 3.9.3 *Если в системе линейных алгебраических уравнений $Ax = b$ квадратная матрица A имеет диагональное преобладание, то метод Якоби для решения этой системы сходится при любом начальном приближении.*

Доказательство. Диагональное преобладание в матрице $A = (a_{ij})$ означает, что

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad i = 1, 2, \dots, n.$$

Следовательно,

$$\sum_{j \neq i} \left| \frac{a_{ij}}{a_{ii}} \right| < 1, \quad i = 1, 2, \dots, n,$$

что равносильно

$$\max_{1 \leq i \leq n} \left(\sum_{j \neq i} \left| \frac{a_{ij}}{a_{ii}} \right| \right) < 1.$$

В выражении, стоящем в левой части неравенства, легко угадать подчинённую чебышёвскую норму (∞ -норму) матрицы $D^{-1}(\tilde{L} + \tilde{U})$, которая была выписана нами в (3.101). Таким образом,

$$\|D^{-1}(\tilde{L} + \tilde{U})\|_{\infty} < 1,$$

откуда, ввиду результата Предложения 3.9.1, следует доказываемое. ■

Итерационный метод Якоби был изобретён в середине XIX века и сейчас при практическом решении систем линейных алгебраических уравнений используется редко, так как существенно проигрывает по эффективности более современным численным методам.²¹ Тем не менее, совсем сбрасывать метод Якоби со счёта будет преждевременным. Лежащая в его основе идея выделения из оператора системы уравнений «диагональной части» достаточно плодотворна и может быть с успехом применена в различных ситуациях.

Рассмотрим, к примеру, систему уравнений

$$Ax = b(x),$$

в которой A — $n \times n$ -матрица, $b(x)$ — некоторая вектор-функция от неизвестной переменной x . В случае, когда $b(x)$ — нелинейная функция, никакие численные методы для решения СЛАУ здесь уже неприменимы, но для отыскания решения мы можем воспользоваться незначительной модификацией итераций Якоби

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

$k = 0, 1, 2, \dots$, с некоторым начальным приближением $x^{(0)}$. Если $b(x)$ изменяется «достаточно медленно», так что $|b'_i(x)/a_{ii}| < 1$ для любых $x \in \mathbb{R}^n$ при всех $i = 1, 2, \dots, n$, то сходимость этого процесса для произвольного начального приближения следует, к примеру, из теоремы Шрёдера о неподвижной точке (Теорема 4.4.5, стр. 463).

²¹Примеры применения и детальные оценки скорости сходимости метода Якоби для решения модельных задач математической физики можно увидеть в [37].

Вообще, *нелинейный итерационный процесс Якоби* в применении к системе уравнений

$$\begin{cases} F_1(x_1, x_2, \dots, x_n) = 0, \\ F_2(x_1, x_2, \dots, x_n) = 0, \\ \vdots \quad \ddots \quad \vdots \\ F_n(x_1, x_2, \dots, x_n) = 0 \end{cases}$$

может заключаться в следующем. Задавшись каким-то начальным приближением $x^{(0)}$, на очередном k -ом шаге для всех $i = 1, 2, \dots, n$ последовательно находят решения \tilde{x}_i уравнений

$$F_i(x_1^{(k)}, \dots, x_{i-1}^{(k)}, x_i, x_{i+1}^{(k)}, \dots, x_n^{(k)}) = 0$$

относительно x_i , а затем полагают $x_i^{(k+1)} \leftarrow \tilde{x}_i$, $i = 1, 2, \dots, n$.

3.9е Итерационный метод Гаусса-Зейделя

В итерационном методе Якоби при организации вычислений по инструкции

$$x_i^{(k+1)} \leftarrow \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij} x_j^{(k)} \right), \quad i = 1, 2, \dots, n, \quad (3.102)$$

компоненты очередного приближения $x^{(k+1)}$ находятся последовательно одна за другой, так что к моменту вычисления i -ой компоненты вектора $x^{(k+1)}$ уже найдены $x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_{i-1}^{(k+1)}$. Но метод Якоби никак не использует эти новые значения, и при вычислении любой компоненты следующего приближения всегда опирается только на вектор $x^{(k)}$ предшествующего приближения. Если итерации сходятся к решению, то естественно ожидать, что все компоненты $x^{(k+1)}$ ближе к искомому решению, чем $x^{(k)}$, а посему немедленное вовлечение их в процесс вычислений будет способствовать ускорению сходимости.

На этой идее основан *итерационный метод Гаусса-Зейделя*,²² псевдокод которого представлен в Табл. 3.6 (где, как и ранее, k — счётчик

²²В отечественной литературе по вычислительной математике нередко используется также термин «метод Зейделя».

Таблица 3.6. Итерационный метод Гаусса-Зейделя
для решения линейных систем уравнений

```

k ← 0;
выбираем начальное приближение x(0);
DO WHILE ( метод не сошёлся )
  DO FOR i = 1 TO n
    xi(k+1) ←  $\frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right)$ 
  END DO
  k ← k + 1;
END DO

```

итераций). В нём суммирование в формуле (3.102) для вычисления i -ой компоненты очередного приближения $x^{(k+1)}$ разбито на две части — по индексам, предшествующим i , и по индексам, следующим за i . Первая часть суммы использует новые вычисленные значения $x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}$, тогда как вторая — компоненты $x_{i+1}^{(k)}, \dots, x_n^{(k)}$ из старого приближения. Метод Гаусса-Зейделя иногда называют также итерационным методом «последовательных смещений», а его основная идея — немедленно вовлекать уже полученную информацию в вычислительный процесс — с успехом применима и для нелинейных итерационных схем.

Чтобы получить для метода Гаусса-Зейделя матричное представление, перепишем его расчётные формулы в виде

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

Используя введённые в §3.9д матрицы \tilde{L} , D и \tilde{U} , на которые разлагается A , можем записать эти формулы в виде

$$(D + \tilde{L})x^{(k+1)} = -\tilde{U}x^{(k)} + b,$$

т. е.

$$x^{(k+1)} = -(D + \tilde{L})^{-1}\tilde{U}x^{(k)} + (D + \tilde{L})^{-1}b, \quad k = 0, 1, 2, \dots \quad (3.103)$$

Таким образом, метод Гаусса-Зейделя можно рассматривать как итерационный метод, порождённый таким расщеплением матрицы СЛАУ в виде $A = G - H$, что $G = D + \tilde{L}$, $H = -\tilde{U}$.

В силу Теоремы 3.9.1 необходимым и достаточным условием сходимости метода Гаусса-Зейделя из любого начального приближения является неравенство

$$\rho((D + \tilde{L})^{-1}\tilde{U}) < 1.$$

Предложение 3.9.4 *Если в системе линейных алгебраических уравнений $Ax = b$ матрица A имеет диагональное преобладание, то метод Гаусса-Зейделя для решения этой системы сходится при любом начальном приближении.*

Доказательство. Отметим, прежде всего, что в условиях диагонального преобладания в A решение x^* рассматриваемой линейной системы существует (вспомним признак неособенности Адамара, §3.2e). Пусть, как и ранее, $x^{(k)}$ — приближение к решению, полученное на k -ом шаге итерационного процесса. Исследуем поведение погрешности решения $z^{(k)} = x^{(k)} - x^*$ в зависимости от номера итерации k .

Чтобы получить формулу для $z^{(k)}$, предварительно перепишем соотношения, которым удовлетворяет точное решение x^* : вместо

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

можно придать им следующий эквивалентный вид

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

Вычитая затем почленно эти равенства из расчётных формул метода Гаусса-Зейделя, т. е. из

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

получим

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

Беря абсолютное значение от обеих частей этого равенства и пользуясь неравенством треугольника для оценки сумм в правой части, будем иметь для $i = 1, 2, \dots, n$:

$$\begin{aligned} |z_i^{(k+1)}| &\leq \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right| \cdot |z_j^{(k+1)}| + \sum_{j=i+1}^n \left| \frac{a_{ij}}{a_{ii}} \right| \cdot |z_j^{(k)}| \\ &\leq \|z^{(k+1)}\|_\infty \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right| + \|z^{(k)}\|_\infty \sum_{j=i+1}^n \left| \frac{a_{ij}}{a_{ii}} \right|. \end{aligned} \quad (3.104)$$

С другой стороны, условие диагонального преобладания в матрице A , т. е.

$$\sum_{j \neq i} |a_{ij}| < |a_{ii}|, \quad i = 1, 2, \dots, n,$$

означает существование константы \varkappa , $0 \leq \varkappa < 1$, такой что

$$\sum_{j \neq i} |a_{ij}| \leq \varkappa |a_{ii}|, \quad i = 1, 2, \dots, n. \quad (3.105)$$

По этой причине

$$\sum_{j \neq i} \left| \frac{a_{ij}}{a_{ii}} \right| \leq \varkappa, \quad i = 1, 2, \dots, n,$$

откуда для $i = 1, 2, \dots, n$ следует

$$\sum_{j=i+1}^n \left| \frac{a_{ij}}{a_{ii}} \right| \leq \varkappa - \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right| \leq \varkappa - \varkappa \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right| = \varkappa \left(1 - \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right| \right).$$

Подставляя полученную оценку в неравенства (3.104), приходим к соотношениям

$$|z_i^{(k+1)}| \leq \|z^{(k+1)}\|_\infty \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right| + \varkappa \|z^{(k)}\|_\infty \left(1 - \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right| \right), \quad (3.106)$$

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

Предположим, что $\max_{1 \leq i \leq n} |z_i^{(k+1)}|$ достигается при $i = l$, так что

$$\|z^{(k+1)}\|_\infty = |z_l^{(k+1)}|. \quad (3.107)$$

Рассмотрим теперь отдельно l -ое неравенство из (3.106). Привлекая равенство (3.107), можем утверждать, что

$$\|z^{(k+1)}\|_\infty \leq \|z^{(k+1)}\|_\infty \sum_{j=1}^{l-1} \left| \frac{a_{lj}}{a_{ll}} \right| + \varkappa \|z^{(k)}\|_\infty \left(1 - \sum_{j=1}^{l-1} \left| \frac{a_{lj}}{a_{ll}} \right| \right),$$

то есть

$$\|z^{(k+1)}\|_\infty \left(1 - \sum_{j=1}^{l-1} \left| \frac{a_{lj}}{a_{ll}} \right| \right) \leq \varkappa \|z^{(k)}\|_\infty \left(1 - \sum_{j=1}^{l-1} \left| \frac{a_{lj}}{a_{ll}} \right| \right). \quad (3.108)$$

Конечно, значение индекса l , на котором достигается равенство (3.107), может меняться в зависимости от номера итерации k . Но так как вплоть до оценки (3.106) мы отслеживали все компоненты погрешности $z_i^{(k+1)}$, то вне зависимости от k неравенство (3.108) должно быть справедливым для компоненты с номером l , определяемой условием (3.107).

Далее, в силу диагонального преобладания в матрице A

$$1 - \sum_{j=1}^{l-1} \left| \frac{a_{lj}}{a_{ll}} \right| > 0,$$

и на эту положительную величину можно сократить обе части неравенства (3.108). Окончательно получаем

$$\|z^{(k+1)}\|_\infty \leq \varkappa \|z^{(k)}\|_\infty,$$

что при $|\varkappa| < 1$ означает сходимость метода Гаусса-Зейделя. ■

Фактически, в доказательстве Предложения 3.9.4 мы получили даже оценку уменьшения чебышёвской нормы погрешности через «меру диагонального преобладания» в матрице СЛАУ, в качестве которой может выступать величина \varkappa , определённая посредством (3.105).

Теорема 3.9.2 *Если в системе линейных алгебраических уравнений $Ax = b$ матрица A является симметричной положительно определённой, то метод Гаусса-Зейделя сходится к решению из любого начального приближения.*

Доказательство может быть найдено, к примеру, в [3, 11]. Теорема 3.9.2 является частным случаем теоремы Островского-Райха (теорема 3.9.3), которая, в свою очередь, может быть получена как следствие из более общей теории итерационных методов, развитой А.А. Самарским и начала которой мы излагаем в §3.12.

Метод Гаусса-Зейделя был сконструирован как модификация метода Якоби, и, казалось бы, должен работать лучше. Так оно и есть «в среднем», на случайно выбранных системах — метод Гаусса-Зейделя работает несколько быстрее, что можно показать математически строго при определённых допущениях на систему. Но в целом ситуация не столь однозначна. Для СЛАУ размера 3×3 и более существуют примеры, на которых метод Якоби расходится, но метод Гаусса-Зейделя сходится, так же как существуют и примеры другого свойства, когда метод Якоби сходится, а метод Гаусса-Зейделя расходится. В частности, для метода Якоби неверна Теорема 3.9.2, и он может расходиться для систем линейных уравнений с симметричными положительно-определёнными матрицами.

По поводу практического применения метода Гаусса-Зейделя можно сказать почти то же самое, что и о методе Якоби в §3.9д. Для решения систем линейных алгебраических уравнений он используется в настоящее время нечасто, но его идея не утратила своего значения и успешно применяется при построении различных итерационных процессов для решения линейных и нелинейных систем уравнений.

3.9ж Методы релаксации

Одним из принципов, который кладётся в основу итерационных методов решения систем уравнений, является так называемый *принцип релаксации*.²³ Он понимается как специальная организация итераций, при которой на каждом шаге процесса уменьшается какая-либо величина, характеризующая погрешность решения системы.

Поскольку само решение x^* нам неизвестно, то оценить напрямую погрешность $(x^{(k)} - x^*)$ не представляется возможным. По этой причине о степени близости $x^{(k)}$ к x^* судят на основании косвенных признаков, важнейшим среди которых является величина *невязки* решения. Невязка определяется как разность левой и правой частей уравнения после подстановки в него приближения к решению, и в нашем случае

²³От латинского слова «relaxatio» — уменьшение напряжения, ослабление.

это $Ax^{(k)} - b$. При этом конкретное применение принципа релаксации может заключаться в том, что на каждом шаге итерационного процесса стремятся уменьшить абсолютные значения компонент вектора невязки либо её норму, либо какую-то зависящую от них величину. В этом смысле методы Якоби и Гаусса-Зейделя можно рассматривать как итерационные процессы, в которых также осуществляется релаксация, поскольку на каждом их шаге компоненты очередного приближения вычисляются из условия зануления соответствующих компонент невязки на основе уже полученной информации о решении. Правда, это делается «локально», для отдельно взятой компоненты, и без учёта влияния результатов вычисления этой компоненты на другие компоненты невязки.

Различают релаксацию *полную* и *неполную*, в зависимости от того, добиваемся ли мы на каждом отдельном шаге итерационного процесса (или его подшаге) наибольшего возможного улучшения рассматриваемой функции от погрешности или нет. Локально полная релаксация может казаться наиболее выгодной, но глобально, с точки зрения сходимости процесса в целом, тщательно подобранная неполная релаксация нередко приводит к более эффективным методам.

Очень популярной реализацией высказанных выше общих идей является метод решения систем линейных алгебраических уравнений, в котором для улучшения сходимости берётся «взвешенное среднее» значений компонент предшествующей $x^{(k)}$ и последующей $x^{(k+1)}$ итераций метода Гаусса-Зейделя. Более точно, зададимся вещественным числом ω , которое будем называть *параметром релаксации*, и i -ую компоненту очередного $(k + 1)$ -го приближения положим равной

$$\omega x_i^{(k+1)} + (1 - \omega)x_i^{(k)},$$

где $x_i^{(k)}$ — i -ая компонента приближения, полученного в результате k -го шага алгоритма, а $x_i^{(k+1)}$ — i -ая компонента приближения, которое было бы получено на основе $x^{(k)}$ и $x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_{i+1}^{(k)}, \dots, x_n^{(k)}$ с помощью метода Гаусса-Зейделя. Псевдокод получающегося итерационного алгоритма, который обычно и называют методом релаксации для решения систем линейных алгебраических уравнений, представлен в Табл. 3.7.

Таблица 3.7. Псевдокод метода релаксации
для решения систем линейных уравнений

```

k ← 0;
выбираем начальное приближение x(0);
DO WHILE ( метод не сошёлся )
  DO FOR i = 1 TO n
    xi(k+1) ← (1 - ω) xi(k)
    +  $\frac{\omega}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right)$ 
  END DO
  k ← k + 1;
END DO

```

Расчётные формулы этого метода можно записать в виде

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

для $i = 1, 2, \dots, n$. Используя введённые выше в §3.9е матрицы \tilde{L} , D и \tilde{U} , можно придать этим соотношениям более компактный вид

$$(D + \omega \tilde{L}) x^{(k+1)} = ((1 - \omega)D - \omega \tilde{U}) x^{(k)} + \omega b,$$

откуда

$$x^{(k+1)} = (D + \omega \tilde{L})^{-1} ((1 - \omega)D - \omega \tilde{U}) x^{(k)} + (D + \omega \tilde{L})^{-1} \omega b,$$

$k = 0, 1, 2, \dots$

В зависимости от конкретного значения параметра релаксации принято различать три случая:

если $\omega < 1$, то говорят о «нижней релаксации»,

если $\omega = 1$, то имеем итерации Гаусса-Зейделя,

если $\omega > 1$, то говорят о «верхней релаксации».²⁴

Последний случай может показаться экзотичным, но во многих ситуациях он действительно обеспечивает улучшение сходимости итераций в сравнении с методом Гаусса-Зейделя. Несколько упрощённое объяснение этого явления может состоять в том, что если направление от $x^{(k)}$ к $x^{(k+1)}$ оказывается удачным в том смысле, что приближает к искомому решению, то имеет смысл пройти по нему и дальше, за $x^{(k+1)}$. Это соответствует случаю $\omega > 1$.

Важно отметить, что метод релаксации также укладывается в изложенную в §3.9в схему итерационных процессов, порождаемых расщеплением матрицы решаемой системы уравнений. Именно, мы берём $A = G_\omega - H_\omega$ с матрицами

$$G_\omega = D + \omega\tilde{L}, \quad H_\omega = (1 - \omega)D - \omega\tilde{U}.$$

Необходимое и достаточное условие сходимости метода релаксации принимает поэтому вид

$$\rho(G_\omega^{-1}H_\omega) < 1.$$

Для некоторых специфичных, но очень важных задач математической физики значение релаксационного параметра ω , при котором величина $\rho(G_\omega^{-1}H_\omega)$ достигает минимума, находится относительно просто. В более сложных задачах для оптимизации ω требуется весьма трудный анализ спектра матрицы перехода $G_\omega^{-1}H_\omega$ из представления (3.97). Обзоры состояния дел в этой области читатель может найти в [45, 49, 77, 95, 96].

Предложение 3.9.5 *Если $C_\omega = (D + \omega\tilde{L})^{-1}((1 - \omega)D - \omega\tilde{U})$ — матрица оператора перехода метода релаксации, то $\rho(C_\omega) \geq |\omega - 1|$. Как следствие, неравенство $0 < \omega < 2$ на параметр релаксации необходимо для сходимости метода.*

Доказательство. Прежде всего, преобразуем матрицу C_ω для прида-

²⁴В англоязычной литературе по вычислительной линейной алгебре этот метод обычно обозначают аббревиатурой SOR(ω), которая происходит от термина «Successive OverRelaxation» — последовательная верхняя релаксация.

ния ей более удобного для дальнейших выкладок вида:

$$\begin{aligned} C_\omega &= (D + \omega\tilde{L})^{-1}((1 - \omega)D - \omega\tilde{U}) \\ &= (D(I + \omega D^{-1}\tilde{L}))^{-1}((1 - \omega)D - \omega\tilde{U}) \\ &= (I + \omega D^{-1}\tilde{L})^{-1}D^{-1}((1 - \omega)D - \omega\tilde{U}) \\ &= (I + \omega D^{-1}\tilde{L})^{-1}((1 - \omega)I - \omega D^{-1}\tilde{U}). \end{aligned}$$

Желая исследовать расположение собственных чисел $\lambda_i(C_\omega)$ матрицы C_ω , рассмотрим её характеристический полином

$$\begin{aligned} \phi(\lambda) &= \det(C_\omega - \lambda I) = \det\left((I + \omega D^{-1}\tilde{L})^{-1}((1 - \omega)I - \omega D^{-1}\tilde{U}) - \lambda I\right) \\ &= p_n \lambda^n + p_{n-1} \lambda^{n-1} + \dots + p_1 \lambda + p_0, \end{aligned}$$

в котором $p_n = (-1)^n$ по построению. Свободный член p_0 характеристического полинома может быть найден как $\phi(0)$:

$$\begin{aligned} p_0 &= \det C_\omega = \det\left((I + \omega D^{-1}\tilde{L})^{-1}((1 - \omega)I - \omega D^{-1}\tilde{U})\right) \\ &= \det((I + \omega D^{-1}\tilde{L})^{-1}) \cdot \det((1 - \omega)I - \omega D^{-1}\tilde{U}) \\ &= \det((1 - \omega)I - \omega D^{-1}\tilde{U}) = (1 - \omega)^n, \end{aligned}$$

коль скоро матрица $(I + \omega D^{-1}\tilde{L})$ — нижняя треугольная и диагональными элементами имеет единицы, а $((1 - \omega)I - \omega D^{-1}\tilde{U})$ — верхняя треугольная, с элементами $(1 - \omega)$ по главной диагонали.

С другой стороны, по теореме Виета свободный член характеристического полинома матрицы, делённый на старший коэффициент, равен произведению его корней, т.е. собственных чисел матрицы, умноженному на $(-1)^n$ (см., к примеру, [22]), и поэтому

$$\prod_{i=1}^n \lambda_i(C_\omega) = (1 - \omega)^n.$$

Отсюда необходимо следует

$$\max_{1 \leq i \leq n} |\lambda_i(C_\omega)| \geq |\omega - 1|,$$

что и доказывает Предложение. ■

Теорема 3.9.3 (теорема Островского-Райха) *Если в системе линейных алгебраических уравнений $Ax = b$ матрица A является симметричной положительно определённой, то для всех значений параметра $\omega \in]0, 2[$ метод релаксации сходится к решению из любого начального приближения.*

Доказательство опускается. Читатель может найти его, к примеру, в книгах [13, 95, 96]. Обоснование теоремы Островского-Райха будет также дано ниже в §3.12 как следствие теоремы Самарского, дающей достаточные условия сходимости для итерационных методов весьма общего вида.

3.10 Нестационарные итерационные методы для линейных систем

3.10a Теоретическое введение

В этом параграфе для решения систем линейных алгебраических уравнений мы рассмотрим нестационарные итерационные методы, которые распространены не меньше стационарных. В основу нестационарных методов могут быть положены различные идеи.

В качестве первого примера рассмотрим метод простой итерации (3.98)

$$x^{(k+1)} \leftarrow (I - \tau A)x^{(k)} + \tau b, \quad k = 0, 1, 2, \dots,$$

исследованный нами в §3.9г. Если переписать его в виде

$$x^{(k+1)} \leftarrow x^{(k)} - \tau(Ax^{(k)} - b), \quad k = 0, 1, 2, \dots, \quad (3.109)$$

то расчёт каждой последующей итерации $x^{(k+1)}$ может трактоваться как вычитание из $x^{(k)}$ поправки, пропорциональной вектору текущей невязки $(Ax^{(k)} - b)$. Но при таком взгляде на итерационный процесс можно попытаться изменять параметр τ в зависимости от шага, т. е. взять $\tau = \tau_k$ переменным, рассмотрев итерации

$$x^{(k+1)} \leftarrow x^{(k)} - \tau_k(Ax^{(k)} - b), \quad k = 0, 1, 2, \dots \quad (3.110)$$

Эту нестационарную версию метода простой итерации часто связывают с именем Л.Ф. Ричардсона, предложившего её идею ещё в 1910 году. Он, к сожалению, не смог развить удовлетворительной теории выбора параметров τ_k , и для решения этого вопроса потребовалось ещё

несколько десятилетий развития вычислительной математики. Отметим, что задача об оптимальном выборе параметров τ_k на группе из нескольких шагов приводит к так называемым чебышёвским циклическим итерационным методам (см. [37, 45, 77]).

Можно пойти по намеченному выше пути дальше, рассмотрев нестационарное обобщение итерационного процесса

$$x^{(k+1)} \leftarrow (I - \Lambda A)x^{(k)} + \Lambda b, \quad k = 0, 1, 2, \dots,$$

который получен в результате матричного предобуславливания исходной системы линейных алгебраических уравнений. Переписав его вычислительную схему в виде

$$x^{(k+1)} \leftarrow x^{(k)} - \Lambda(Ax^{(k)} - b), \quad k = 0, 1, 2, \dots,$$

нетрудно увидеть возможность изменения предобуславливающей матрицы Λ в зависимости от номера шага. Таким образом, приходим к весьма общей схеме нестационарных линейных итерационных процессов

$$x^{(k+1)} \leftarrow x^{(k)} - \Lambda_k(Ax^{(k)} - b), \quad k = 0, 1, 2, \dots,$$

где $\{\Lambda_k\}_{k=0}^{\infty}$ — некоторая последовательность матриц, выбор которой зависит, вообще говоря, от начального приближения $x^{(0)}$.

Другой популярный путь построения нестационарных итерационных методов для решения уравнений — использование *вариационных принципов*.

Интуитивно понятный термин «вариация» был введён в математику Ж.-Л. Лагранжем для обозначения малого изменения («шевеления») независимой переменной или рассматриваемой функции (функционала). Соответственно, метод исследования задач нахождения экстремумов, основанный на изучении зависимости функции от вариаций аргументов получил название *метода вариаций*. Но со временем «вариационными» стали именовать методы решения различных уравнений, которые сводят исходную постановку задачи к определённым задачам нахождение экстремума. Согласно этой терминологии, *вариационными принципами* теперь называют переформулировки интересующих нас задач в виде каких-либо оптимизационных задач, т. е. задач нахождение минимумов или максимумов. Тогда итерационные методы решения СЛАУ могут конструироваться как итерационные процессы для отыскания этих экстремумов тех или иных функционалов.

Вариационные принципы получаются весьма различными способами. Некоторые из них вытекают из содержательного (физического, механического и пр.) смысла решаемой задачи. Например, в классической механике хорошо известны «принцип наименьшего действия Лагранжа», в оптике существует «принцип Ферма» [70]. В последнее столетие имеется тенденция всё меньше связывать вариационные принципы с конкретным физическим содержанием, они становятся абстрактным математическим инструментом решения разнообразных задач.

Строго говоря, в вычислительном отношении получающаяся в результате описанного выше сведения оптимизационная задача может быть не вполне эквивалентна исходной, так как задача нахождения устойчивого решения уравнения может превратиться в неустойчивую задачу о проверке точного равенства экстремума нулю (этот вопрос более подробно обсуждается далее в §4.2б). Но если существование решения уравнения известно априори, до того, как мы приступаем к его нахождению (например, на основе каких-либо теорем существования), то вариационные методы становятся важным подспорьем практических вычислений. Именно такова ситуация с системами линейных алгебраических уравнений, разрешимость которых часто обеспечивается различными результатами из линейной алгебры.

Как именно можно переформулировать задачу решения СЛАУ в виде оптимизационной задачи? По-видимому, простейший способ может основываться на том факте, что точное решение x^* зануляет норму невязки $\|Ax - b\|$, доставляя ей, таким образом, наименьшее возможное значение. Желая приобрести гладкость получаемого функционала по неизвестной переменной x , обычно берут евклидову норму невязки, или даже её квадрат, т. е. скалярное произведение $\langle Ax - b, Ax - b \rangle$, чтобы не привлекать взятия корня. Получающаяся задача минимизации величины $\|Ax - b\|_2^2$ является называется *линейной задачей о наименьших квадратах*, и мы рассмотрим её подробнее в §3.15.

Ещё одним фактом, который служит теоретической основой для вариационных методов решения систем линейных алгебраических уравнений является

Предложение 3.10.1 *Вектор $x^* \in \mathbb{R}^n$ является решением системы линейных алгебраических уравнений $Ax = b$ с симметричной положительно определённой матрицей A тогда и только тогда, когда он доставляет минимум функционалу $\Phi(x) = \frac{1}{2}\langle Ax, x \rangle - \langle b, x \rangle$.*

Доказательство. Если A — симметричная положительно-определённая матрица, то, как мы видели в §3.3а, выражением $\frac{1}{2}\langle Ax, x \rangle$ задаётся так называемая энергетическая норма $\|\cdot\|_A$ векторов из \mathbb{R}^n .

Далее, пусть x^* — решение рассматриваемой системы линейных алгебраических уравнений $Ax = b$, которое существует и единственно в силу положительной определённости матрицы A . Из единственности x^* следует, что некоторый вектор $x \in \mathbb{R}^n$ является решением системы уравнений тогда и только тогда, когда $x - x^* = 0$, или, иными словами, $\|x - x^*\|_A^2 = 0$.

С другой стороны, учитывая симметричность матрицы A и равенство $Ax^* = b$, получаем

$$\begin{aligned} \Phi(x) &= \frac{1}{2}\langle Ax, x \rangle - \langle b, x \rangle \\ &= \frac{1}{2}\langle Ax, x \rangle - \langle Ax^*, x \rangle + \frac{1}{2}\langle Ax^*, x^* \rangle - \frac{1}{2}\langle Ax^*, x^* \rangle \\ &= \frac{1}{2}(\langle Ax, x \rangle - \langle Ax, x^* \rangle - \langle Ax^*, x \rangle + \langle Ax^*, x^* \rangle) - \frac{1}{2}\langle Ax^*, x^* \rangle \\ &= \frac{1}{2}\langle A(x - x^*), x - x^* \rangle - \frac{1}{2}\langle Ax^*, x^* \rangle \\ &= \frac{1}{2}\|x - x^*\|_A^2 - \frac{1}{2}\langle Ax^*, x^* \rangle, \end{aligned} \tag{3.111}$$

так что функционал $\Phi(x)$ отличается от половины квадрата энергетической нормы погрешности лишь постоянным слагаемым $\frac{1}{2}\langle Ax^*, x^* \rangle$ (которое заранее неизвестно из-за незнания нами x^*). Следовательно, $\Phi(x)$ действительно достигает своего единственного минимума при том же значении аргумента, что и $\|x - x^*\|_A^2$, т. е. на решении x^* рассматриваемой линейной системы. ■

Функционал $\Phi(x) = \frac{1}{2}\langle Ax, x \rangle - \langle b, x \rangle$, который является квадратичной формой от вектора переменных x , обычно называют *функционалом энергии* из-за его сходства с выражениями для различных видов энергии в физических системах. К примеру, кинетическая энергия тела массы m , движущегося со скоростью v , равна $\frac{1}{2}mv^2$. Энергия упругой деформации пружины с жёсткостью k , растянутой или сжатой на величину x , равна $\frac{1}{2}kx^2$, и т. п.

Поскольку A — симметричная матрица, то ортогональным преобразованием подобия она может быть приведена к диагональной матрице D , на главной диагонали которой стоят собственные значения λ :

$$A = Q^T D Q,$$

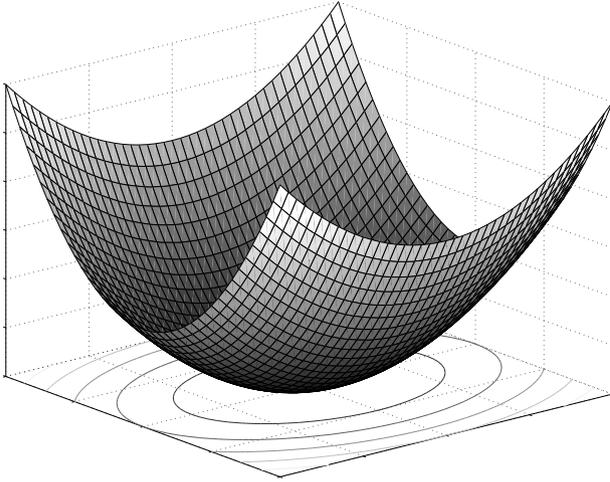


Рис. 3.20. Типичный график функционала энергии и его линии уровня.

причём в силу положительной определённости матрицы A диагональные элементы D положительны. Подставляя это представление в выражение для функционала энергии $\Phi(x)$, получим

$$\begin{aligned}\Phi(x) &= \langle Q^{\top} D Q x, x \rangle - 2\langle b, x \rangle \\ &= \langle D(Qx), Qx \rangle - 2\langle Qb, Qx \rangle \\ &= \langle Dy, y \rangle - 2\langle Qb, y \rangle,\end{aligned}$$

где обозначено $y = Qx$. Видим, что в изменённой системе координат, которая получается с помощью ортогонального линейного преобразования переменных, выражение для функционала энергии $\Phi(x)$ есть сумма квадратов с коэффициентами, равными собственным значениям матрицы A , т. е. член $\langle Dy, y \rangle$, минус линейный член $2\langle Qb, y \rangle$. Таким образом, график функционала энергии — это эллиптический параболоид, возможно, сдвинутый относительно начала координат и ещё повернутый, а его поверхности уровня (линии уровня в двумерном случае) — эллипсоиды (эллипсы), в центре которых находится искомое решение системы уравнений. При этом форма эллипсоидов уровня находится в зависимости от разброса коэффициентов при квадратах переменных,

т.е. от числа обусловленности матрицы A . Чем больше эта обусловленность, тем сильнее сплющены эллипсоиды уровня, так что для плохообусловленных СЛАУ решение находится на дне длинного и узкого «оврага».

3.106 Метод наискорейшего спуска

В предшествующем пункте были предложены две вариационные переформулировки задачи решения системы линейных алгебраических уравнений. Как находить минимум соответствующих функционалов? Прежде, чем строить конкретные численные алгоритмы, рассмотрим общую схему.

Пусть $f : \mathbb{R}^n \rightarrow \mathbb{R}$ — некоторая функция, ограниченная снизу на всём пространстве \mathbb{R}^n и принимающая своё наименьшее значение в x^* , так что

$$f(x) \geq f(x^*) = \min_{x \in \mathbb{R}^n} f(x) \quad \text{для любых } x \in \mathbb{R}^n.$$

Нам нужно найти точку x^* . При этом саму функцию f , для которой ищется экстремум, в теории оптимизации называют *целевой функцией*.

Различают экстремумы *локальные* и *глобальные*. Локальными называют экстремумы, в которых значения целевой функции лучше, чем в некоторой окрестности рассматриваемой точки. Глобальные экстремумы доставляют функции значения, лучшие всех значений функции на всей её области определения. Нас в связи с задачей минимизации функционала энергии интересуют, конечно, его глобальные минимумы.

Типичным подходом к решению задач оптимизации является итерационное построение последовательности значений аргумента $\{x^{(k)}\}$, которая «минимизирует» функцию f в том смысле, что

$$\lim_{k \rightarrow \infty} f(x^{(k)}) = \min_{x \in \mathbb{R}^n} f(x).$$

Если построенная последовательность $\{x^{(k)}\}$ сходится к некоторому пределу, то он и является решением задачи x^* в случае непрерывной функции f .

Метод градиентного спуска, является способом построения последовательности, которая является минимизирующей для определённого класса дифференцируемых целевых функций, и заключается в следующем. Пусть уже найдено какое-то приближение $x^{(k)}$, $k = 0, 1, 2, \dots$, к

точке минимума функции $f(x)$. Естественная идея состоит в том, чтобы из $x^{(k)}$ сдвинуться по направлению наибольшего убывания целевой функции, которое противоположно направлению градиента $f'(x^{(k)})$, т.е. взять

$$x^{(k+1)} \leftarrow x^{(k)} - \tau_k f'(x^{(k)}), \quad (3.112)$$

где τ_k — величина шага, которая выбирается из условия убывания целевой функции на рассматриваемой итерации. Далее мы можем повторить этот шаг ещё раз и ещё ... столько, сколько требуется для достижения требуемого приближения к минимуму.

Если целевая функция имеет более одного локального экстремума, то этот метод может сходиться к какому-нибудь одному из них, который не обязательно является глобальным. К счастью, подобный феномен не может случиться в интересующем нас случае минимизации функционала энергии $\Phi(x)$, порождаемого системой линейных уравнений с симметричной положительно определённой матрицей. Свойства $\Phi(x)$ достаточно хороши, и он имеет один локальный минимум, который одновременно и глобален.

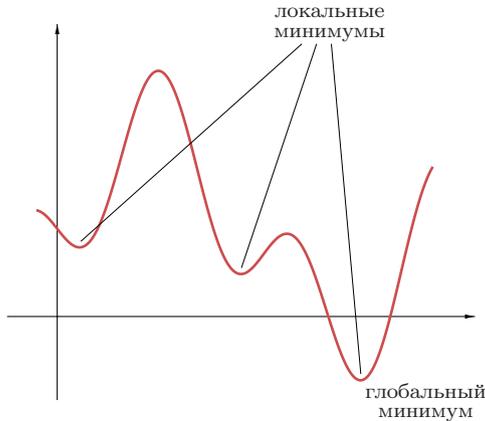


Рис. 3.21. Глобальные и локальные минимумы функции.

Вычислим градиент функционала энергии:

$$\frac{\partial \Phi(x)}{\partial x_i} = \frac{\partial}{\partial x_i} \left(\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j - \sum_{i=1}^n b_i x_i \right) = \sum_{j=1}^n a_{ij} x_j - b_i,$$

$l = 1, 2, \dots, n$. Множитель $1/2$ исчезает в результате потому, что в двойной сумме помимо квадратичных слагаемых $a_{ii}x_i^2$ остальные слагаемые присутствуют парами, как $a_{ij}x_i x_j$ и $a_{ji}x_j x_i$, причём $a_{ij} = a_{ji}$. В целом

$$\Phi'(x) = \left(\frac{\partial \Phi(x)}{\partial x_1}, \frac{\partial \Phi(x)}{\partial x_2}, \dots, \frac{\partial \Phi(x)}{\partial x_n} \right)^\top = Ax - b,$$

т. е. градиент функционала Φ равен невязке решаемой системы линейных уравнений в рассматриваемой точке. Важнейшим выводом из этого факта является тот, что метод простой итерации (3.98)–(3.109) является ни чем иным, как методом градиентного спуска (3.112) для минимизации функционала энергии Φ , в котором шаг τ_k выбран постоянным и равным τ . Вообще, метод градиентного спуска (3.112) оказывается равносильным простейшему нестационарному итерационному методу (3.110).

Выбор величины шага τ_k является очень ответственным делом, так как от него зависит и наличие сходимости, и её скорость. Спуск по направлению антиградиента обеспечивает убывание целевой функции лишь при достаточно малых шагах, и потому при неудачно большой величине шага мы можем попасть в точку, где значение функционала не меньше, чем в текущей точке. С другой стороны, слишком малый шаг приведёт к очень медленному движению в сторону решения. Для градиентного метода с постоянным шагом его трактовка как метода простой итерации позволяет, опираясь на результат §3.9г, выбрать шаг $\tau_k = \text{const}$, который наверняка обеспечивает сходимость процесса. Именно, если положительные числа μ и M — это нижняя и верхняя границы спектра положительно определённой матрицы A решаемой системы, то в соответствии с (3.99) для сходимости следует взять

$$\tau_k = \tau = \frac{2}{M + \mu}.$$

Другой способ выбора шага состоит в том, чтобы потребовать τ_k наибольшим возможным, обеспечивающим убывание функционала Φ вдоль выбранного направления спуска по антиградиенту. При этом получается разновидность градиентного спуска, называемая *методом наискорейшего спуска*, теория которого была разработана в конце 40-х годов XX века Л.В. Канторовичем.

Для определения конкретной величины шага τ_k в методе наискорейшего спуска нужно подставить выражение $x^{(k)} - \tau_k \Phi'(x^{(k)}) = x^{(k)} -$

$\tau_k(Ax^{(k)} - b)$ в аргумент функционала энергии и продифференцировать получающееся отображение по τ_k . Для удобства выкладок обозначим невязку $r^{(k)} := Ax^{(k)} - b$. Имеем

$$\begin{aligned}\Phi(x^{(k)} - \tau_k r^{(k)}) &= \frac{1}{2} \langle Ax^{(k)} - \tau_k r^{(k)}, x^{(k)} - \tau_k r^{(k)} \rangle - \langle b, x^{(k)} - \tau_k r^{(k)} \rangle \\ &= \frac{1}{2} \langle Ax^{(k)}, x^{(k)} \rangle - \tau_k \langle Ax^{(k)}, r^{(k)} \rangle + \frac{1}{2} \tau_k^2 \langle Ar^{(k)}, r^{(k)} \rangle \\ &\quad - \langle b, x^{(k)} \rangle + \tau_k \langle b, r^{(k)} \rangle.\end{aligned}$$

При дифференцировании выписанного выражения по τ_k не зависящие от него члены исчезнут, и мы получим

$$\begin{aligned}\frac{d}{d\tau_k} \Phi(x^{(k)} - \tau_k r^{(k)}) &= -\langle Ax^{(k)}, r^{(k)} \rangle + \tau_k \langle Ar^{(k)}, r^{(k)} \rangle + \langle b, r^{(k)} \rangle \\ &= \tau_k \langle Ar^{(k)}, r^{(k)} \rangle - \langle Ax^{(k)} - b, r^{(k)} \rangle \\ &= \tau_k \langle Ar^{(k)}, r^{(k)} \rangle - \langle r^{(k)}, r^{(k)} \rangle.\end{aligned}$$

Таким образом, в точке экстремума по τ_k из условия

$$\frac{d}{d\tau_k} \Phi(x^{(k)} - \tau_k r^{(k)}) = 0$$

необходимо следует

$$\tau_k = \frac{\langle r^{(k)}, r^{(k)} \rangle}{\langle Ar^{(k)}, r^{(k)} \rangle}.$$

Легко видеть, что при найденном значении τ_k функционалом энергии действительно достигается минимум по выбранному направлению спуска, так как тогда его вторая производная по τ_k , равная $\langle Ar^{(k)}, r^{(k)} \rangle$, положительна. В целом, псевдокод метода наискорейшего градиентного спуска для решения системы линейных алгебраических уравнений $Ax = b$ представлен в Табл. 3.8.

Теорема 3.10.1 *Если A — симметричная положительно определённая матрица, то последовательность $\{x^{(k)}\}$, порождаемая методом наискорейшего спуска, сходится к решению x^* системы уравнений $Ax = b$ из любого начального приближения $x^{(0)}$, и быстрота этой сходимости оценивается неравенством*

$$\|x^{(k)} - x^*\|_A \leq \left(\frac{M - \mu}{M + \mu} \right)^k \|x^{(0)} - x^*\|_A, \quad k = 0, 1, 2, \dots, \quad (3.113)$$

где μ , M — нижняя и верхняя границы спектра матрицы A .

Таблица 3.8. Псевдокод метода наискорейшего спуска для решения систем линейных уравнений

```

k ← 0;
выбираем начальное приближение x(0);
DO WHILE ( метод не сошёлся )
    r(k) ← Ax(k) - b;
    τk ←  $\frac{\|r^{(k)}\|_2^2}{\langle Ar^{(k)}, r^{(k)} \rangle}$ ;
    x(k+1) ← x(k) - τkr(k);
    k ← k + 1;
END DO

```

Доказательство оценки (3.113) и теоремы в целом будет получено путём сравнения метода наискорейшего спуска с методом градиентного спуска с постоянным оптимальным шагом.

Пусть в результате выполнения $(k-1)$ -го шага метода наискорейшего спуска получено приближение $x^{(k)}$, и мы делаем k -ый шаг, который даёт $x^{(k+1)}$. Обозначима также через \tilde{x} результат выполнения с $x^{(k)}$ одного шага метода простой итерации, так что

$$\tilde{x} = x^{(k)} - \tau(Ax^{(k)} - b).$$

Из развитой в предшествующей части параграфа теории вытекает, что при любом выборе параметра τ

$$\Phi(x^{(k+1)}) \leq \Phi(\tilde{x}).$$

Далее, из равенства (3.111)

$$\Phi(x) = \frac{1}{2}\|x - x^*\|_A^2 - \frac{1}{2}\langle Ax^*, x^* \rangle$$

с постоянным вычитаемым $\frac{1}{2}\langle Ax^*, x^* \rangle$ следует, что

$$\frac{1}{2}\|x^{(k+1)} - x^*\|_A^2 \leq \frac{1}{2}\|\tilde{x} - x^*\|_A^2,$$

т. е.

$$\|x^{(k+1)} - x^*\|_A \leq \|\tilde{x} - x^*\|_A. \quad (3.114)$$

Иными словами, метод, обеспечивающий лучшее убывание значения функционала энергии одновременно обеспечивает лучшее приближение к решению в энергетической норме.

В методе градиентного спуска с постоянным шагом — совпадающем с методом простой итерации (3.98) или (3.109) — имеем

$$\tilde{x} - x^* = (I - \tau A)(x^{(k)} - x^*), \quad k = 0, 1, 2, \dots$$

Матрица $(I - \tau A)$ является многочленом первой степени от матрицы A , и потому можем применить неравенство (3.25) из Предложения 3.3.8 (стр. 248):

$$\|\tilde{x} - x^*\|_A \leq \|I - \tau A\|_2 \|x^{(k)} - x^*\|_A.$$

При этом у метода наискорейшего спуска оценка заведомо не хуже этой оценки, в которой взято значение параметра шага $\tau = 2/(M + \mu)$, оптимальное для спуска с постоянным шагом. Тогда в соответствии с (3.114) и с оценкой (3.100) для метода простой итерации получаем

$$\|x^{(k+1)} - x^*\|_A \leq \left(\frac{M - \mu}{M + \mu} \right) \|x^{(k)} - x^*\|_A, \quad k = 0, 1, 2, \dots,$$

откуда следует доказываемая оценка (3.113). ■

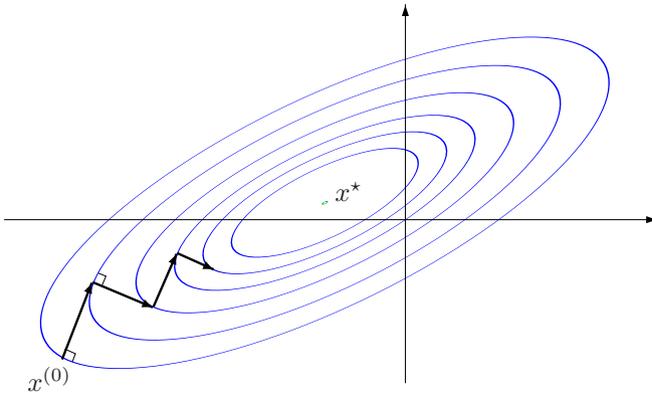


Рис. 3.22. Иллюстрация работы метода наискорейшего спуска.

Интересно и поучительно рассмотреть геометрическую иллюстрацию работы метода наискорейшего спуска.

Градиент функционала энергии нормален к его поверхностям уровня, и именно по этим направлениям осуществляется «спуск» — движение в сторону решения. Шаг в методе наискорейшего спуска идёт на максимально возможную величину — до пересечения с касательным эллипсоидом. Поэтому траектория метода наискорейшего спуска является ломаной, звенья которой перпендикулярны друг другу (см. Рис. 3.22).

Хотя доказательство Теоремы 3.10.1 основано на мажоризации наискорейшего спуска методом простой итерации и может показаться довольно грубым, оценка (3.113) в действительности весьма точно передаёт особенности поведения метода, а именно, замедление сходимости при $M \gg \mu$. Тот факт, что в случае плохой обусловленности матрицы системы движение к решению в методе наискорейшего спуска весьма далеко от оптимального, подтверждается вычислительной практикой и может быть понято на основе геометрической интерпретации. Искомое решение находится при этом на дне глубокого и вытянутого оврага, а метод «рыскает» от одного склона оврага к другому вместо того, чтобы идти напрямую к глубочайшей точке — решению.

3.10в Метод минимальных невязок

Другой популярный подход к выбору итерационных параметров τ_k в нестационарном итерационном процессе (3.110)

$$x^{(k+1)} \leftarrow x^{(k)} - \tau_k (Ax^{(k)} - b), \quad k = 0, 1, 2, \dots,$$

был предложен С.Г. Крейном и М.А. Красносельским в работе [24] и назван ими *методом минимальных невязок*. Его псевдокод приведён в Табл. 3.9. Каждый шаг этого метода минимизирует $\|Ax - b\|_2$ или, что равносильно, $\|Ax - b\|_2^2$ в направлении невязки k -го приближения, равной $r^{(k)} = Ax^{(k)} - b$. Оказывается, что это эквивалентно наибольшему возможному уменьшению $A^\top A$ -нормы погрешности приближённого решения системы. В самом деле, если x^* — точное решение системы

уравнений, то $Ax^* = b$, и потому

$$\begin{aligned} \|Ax - b\|_2^2 &= \langle Ax - b, Ax - b \rangle = \langle Ax - Ax^*, Ax - Ax^* \rangle \\ &= \langle A(x - x^*), A(x - x^*) \rangle = \langle A^\top A(x - x^*), x - x^* \rangle \\ &= \|x - x^*\|_{A^\top A}^2. \end{aligned} \quad (3.115)$$

Если уже найдено $x^{(k)}$, и мы желаем выбрать параметр τ так, чтобы на следующем приближении $x^{(k)} - \tau r^{(k)}$ минимизировать 2-норму невязки решения, то необходимо найти минимум по τ для выражения

$$\begin{aligned} \|A(x^{(k)} - \tau r^{(k)}) - b\|_2^2 &= \langle A(x^{(k)} - \tau r^{(k)}) - b, A(x^{(k)} - \tau r^{(k)}) - b \rangle \\ &= \tau^2 \langle Ar^{(k)}, Ar^{(k)} \rangle - 2\tau (\langle Ax^{(k)}, Ar^{(k)} \rangle - \langle b, Ar^{(k)} \rangle) \\ &\quad + \langle Ax^{(k)}, Ax^{(k)} \rangle + \langle b, b \rangle. \end{aligned}$$

Дифференцируя его по τ и приравнивая производную нулю, получим

$$2\tau \langle Ar^{(k)}, Ar^{(k)} \rangle - 2(\langle Ax^{(k)}, Ar^{(k)} \rangle - \langle b, Ar^{(k)} \rangle) = 0,$$

что с учётом равенства $Ax^{(k)} - b = r^{(k)}$ даёт

$$\tau \langle Ar^{(k)}, Ar^{(k)} \rangle - \langle r^{(k)}, Ar^{(k)} \rangle = 0.$$

Окончательно

$$\tau = \frac{\langle Ar^{(k)}, r^{(k)} \rangle}{\langle Ar^{(k)}, Ar^{(k)} \rangle} = \frac{\langle Ar^{(k)}, r^{(k)} \rangle}{\|Ar^{(k)}\|_2^2}.$$

Теорема 3.10.2 Если A — симметричная положительно определённая матрица, то последовательность $\{x^{(k)}\}$, порождаемая методом минимальных невязок, сходится к решению x^* системы уравнений $Ax = b$ из любого начального приближения $x^{(0)}$, и быстрота этой сходимости оценивается неравенством

$$\|x^{(k)} - x^*\|_{A^\top A} \leq \left(1 - \left(\frac{\mu}{M}\right)^2\right)^{k/2} \|x^{(0)} - x^*\|_{A^\top A}, \quad (3.116)$$

$k = 0, 1, 2, \dots$, где μ , M — нижняя и верхняя границы спектра матрицы A .

Таблица 3.9. Псевдокод метода минимальных невязок для решения систем линейных уравнений

```

k ← 0;
выбираем начальное приближение x(0);
DO WHILE ( метод не сошёлся )
    r(k) ← Ax(k) - b;
    τk ←  $\frac{\langle Ar^{(k)}, r^{(k)} \rangle}{\|Ar^{(k)}\|_2^2}$ ;
    x(k+1) ← x(k) - τkr(k);
    k ← k + 1;
END DO

```

Доказательство теоремы можно найти, к примеру, в книге [56], где для невязок $r^{(k)} = Ax^{(k)} - b$ доказывается оценка

$$\|r^{(k+1)}\|_2 \leq \left(1 - \left(\frac{\mu}{M}\right)^2\right)^{1/2} \|r^{(k)}\|_2, \quad k = 0, 1, 2, \dots$$

С учётом выкладок (3.115) этот результат совершенно равносильен неравенству (3.116).

Для систем линейных уравнений с несимметричными матрицами, которые положительно определены, метод минимальных невязок также сходится. Но если матрица системы не является положительно определённой метод может не сходиться к решению.

Пример 3.10.1 В системе линейных алгебраических уравнений

$$\begin{pmatrix} 2 & 2 \\ 1 & 0 \end{pmatrix} x = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

матрица не является ни симметричной, ни положительно определённой (её собственные значения приблизительно равны 2.732 и -0.7321).

В применении к этой системе метод минимальных невязок с нулевым начальным приближением вскоре после начала работы устанавливается на векторе $(0.7530088, 0.2469912)^\top$, тогда как настоящее решение — это вектор $(1, 0)^\top$. Из других начальных приближений метод будет сходиться к другим векторам, которые также не совпадают с этим точным решением. ■

Практически важной особенностью метода минимальных невязок является быстрая сходимость к решению на первых шагах, которая затем замедляется и выходит на асимптотическую скорость, описываемую Теоремой 3.10.2.

Если сходимость методов наискорейшего спуска и минимальных невязок принципиально не лучше сходимости метода простой итерации, то имеют ли они какое-либо практическое значение? Ответ на этот вопрос положителен. Вспомним, что наша оптимизация метода простой итерации основывалась на знании границ спектра симметричной положительно определённой матрицы СЛАУ. Для работы методов наискорейшего спуска и минимальных невязок этой информации не требуется.

Метод минимальных невязок в представленной выше версии не отличается большой эффективностью. Но он послужил основой для создания многих популярных современных методов решения СЛАУ. В частности, большое распространение на практике получила модификация метода минимальных невязок, известная под англоязычной аббревиатурой GMRES — Generalized Minimal RESiduals — обобщённый метод минимальных невязок, предложенная Ю. Саадом [56] (см. также [43, 59]).

3.10г Метод сопряжённых градиентов

Методами сопряжённых направлений для решения систем линейных алгебраических уравнений вида $Ax = b$ называют методы, в которых решение ищется в виде линейной комбинации векторов, ортогональных в скалярном произведении, которое порождено матрицей системы или же какой-либо матрицей, связанной с матрицей системы. Таким образом, при этом

$$x = x^{(0)} + \sum_{i=1}^n c_i s^{(i)},$$

где $x^{(0)}$ — начальное приближение, $s^{(i)}$, $i = 1, 2, \dots, n$, — векторы «сопряжённых направлений», c_i — коэффициенты разложения решения по ним. Термин «сопряжённые направления» имеет происхождение в аналитической геометрии, где направления, задаваемые векторами u и v , называются сопряжёнными относительно поверхности второго порядка, задаваемой уравнением $\langle Rx, x \rangle = \text{const}$ с симметричной матрицей R , если $\langle Ru, v \rangle = 0$. В методах сопряжённых направлений последовательно строится базис из векторов $s^{(i)}$ и одновременно находятся коэффициенты c_i , $i = 1, 2, \dots, n$.

Наиболее популярными представителями методов сопряжённых направлений являются *методы сопряжённых градиентов*, предложенные М.Р. Хестенсом и Э.Л. Штифелем в начале 50-х годов прошлого века. Их общая схема такова.

Пусть требуется найти решение системы линейных алгебраических уравнений

$$Ax = b$$

с симметричной и положительно определённой матрицей A . Для такой матрицы имеет смысл понятие A -ортогональности, и пусть $s^{(1)}$, $s^{(2)}$, \dots , $s^{(n)}$ — базис \mathbb{R}^n , составленный из A -ортогональных векторов. Решение x^* системы уравнений можно искать в виде разложения по этому базису, т. е.

$$x^* = \sum_{i=1}^n x_i s^{(i)} \quad (3.117)$$

с какими-то неизвестными коэффициентами x_i , $i = 1, 2, \dots, n$. Умножая обе части этого равенства слева на матрицу A и учитывая, что $Ax^* = b$, будем иметь

$$\sum_{i=1}^n x_i (As^{(i)}) = b.$$

Если далее умножить скалярно это равенство на $s^{(j)}$, $j = 1, 2, \dots, n$, то получим n штук соотношений

$$\sum_{i=1}^n x_i \langle As^{(i)}, s^{(j)} \rangle = \langle b, s^{(j)} \rangle, \quad j = 1, 2, \dots, n. \quad (3.118)$$

Но в силу A -ортогональности системы векторов $s^{(1)}$, $s^{(2)}$, \dots , $s^{(n)}$

$$\langle As^{(i)}, s^{(j)} \rangle = \langle s^{(i)}, s^{(j)} \rangle_A = \delta_{ij} = \begin{cases} 0, & \text{если } i \neq j, \\ 1, & \text{если } i = j, \end{cases}$$

так что от равенств (3.118) останется лишь

$$x_i \langle As^{(i)}, s^{(i)} \rangle = \langle b, s^{(i)} \rangle, \quad i = 1, 2, \dots, n.$$

Окончательно

$$x_i = \frac{\langle b, s^{(i)} \rangle}{\langle As^{(i)}, s^{(i)} \rangle}, \quad i = 1, 2, \dots, n,$$

откуда из (3.117) нетрудно восстановить искомое решение СЛАУ. Но для практического применения этого элегантного результата нужно уметь эффективно строить A -ортогональный базис $s^{(1)}, s^{(2)}, \dots, s^{(n)}$ пространства \mathbb{R}^n .

Он определяется процессом A -ортогонализации невязок $r^{(0)}, r^{(1)}, \dots, r^{(n-1)}$ последовательных приближений к решению $x^{(0)}, x^{(1)}, \dots, x^{(n-1)}$. Этот процесс ортогонализации конечен и завершается при некотором $k \leq n$, для которого $r^{(k)} = 0$, т. е. когда очередная невязка приближённого решения зануляется.

Но на практике из-за неизбежных погрешностей вычислений метод сопряжённых градиентов может не прийти к решению системы за n шагов. Тогда целесообразно повторить цикл уточнения, превратив алгоритм при необходимости в итерационный. Именно такой псевдокод метода сопряжённых градиентов приведён в Табл. 3.10. В теле цикла первая команда вычисляет длину очередного шага метода, а вторая строка даёт следующее приближение к решению. Далее вычисляется невязка вновь найденного приближённого решения, а в следующих двух строках тела цикла (перед увеличением счётчика k) вычисляется новое направление движения к решению.

Широко распространена также другая трактовка метода сопряжённых градиентов, представляющая его как модификацию метода наискорейшего градиентного спуска. Как мы видели в предшествующем параграфе, направления градиентов энергетического функционала, по которым осуществляется движение (спуск) к решению в методе наискорейшего спуска, могут сильно изменяться от шага к шагу. По этой причине траектория метода наискорейшего спуска имеет зигзагообразный вид, и для получения решения затрачивается много лишней работы. Естественно попытаться каким-нибудь образом сгладить «вихляния» метода наискорейшего спуска, чтобы он шёл к решению более прямым путём. Один из возможных способов сделать это состоит в том, чтобы на каждой итерации корректировать направление спуска по антигради-

Таблица 3.10. Псевдокод метода сопряжённых градиентов для решения систем линейных уравнений

```

k ← 0;
выбираем начальное приближение x(0);
r(0) ← Ax(0) - b;
s(0) ← r(0);
DO WHILE ( метод не сошёлся )
    τk ←  $\frac{\langle r^{(k)}, r^{(k)} \rangle}{\langle As^{(k)}, s^{(k)} \rangle}$ ;
    x(k+1) ← x(k) - τks(k);
    r(k+1) ← r(k) - τkAs(k);
    υk ←  $\frac{\langle r^{(k+1)}, r^{(k+1)} \rangle}{\langle r^{(k)}, r^{(k)} \rangle}$ ;
    s(k+1) ← r(k+1) + υks(k);
    k ← k + 1;
END DO

```

енту с помощью некоторой добавки. Например, исходя из геометрических соображений, её можно взять пропорциональной разности двух последовательных приближений, так что в целом получаем алгоритм

$$x^{(k+1)} \leftarrow x^{(k)} - \tau_k (Ax^{(k)} - b) + \nu_k (x^{(k)} - x^{(k-1)}), \quad (3.119)$$

$$k = 0, 1, 2, \dots,$$

где τ_k , ν_k — некоторые параметры. Для их определения можно привлечь условие минимизации энергетического функционала $\Phi(x)$ в точке $x^{(k+1)}$. При этом получаются формулы для τ_k и ν_k , приведённые в псевдокоде Табл. 3.10.

Итерационный процесс (3.119) — двухшаговый, так что для начала его работы требуется знать два последовательных приближения к решению. Можно положить $x^{(-1)} = x^{(0)}$, откуда однозначно определяется

$x^{(1)}$ и т. д.

3.11 Методы установления

Методы установления — общее название для большой группы методов, в основе которых лежит идея искать решение рассматриваемой стационарной задачи как предела по времени $t \rightarrow \infty$ для решения связанной с ней вспомогательной нестационарной задачи. Этот подход к решению различных задач был развит в 30-е годы XX века А.Н. Тихоновым.

Пусть требуется решить систему уравнений

$$Ax = b.$$

Наряду с ней рассмотрим также систему уравнений

$$\frac{\partial x(t)}{\partial t} + Ax(t) = b, \quad (3.120)$$

в которой вектор неизвестных переменных x зависит от времени t . Ясно, что если $x(t)$ не зависит от переменной t , то производная $\partial x/\partial t$ зануляется и соответствующие значения $x(t)$ являются решением исходной задачи

Наиболее часто задачу (3.120) рассматривают на бесконечном интервале $[t_0, \infty)$ и ищут её устанавливающееся решение, т. е. такое, что существует конечный $\lim_{t \rightarrow \infty} x(t) = x^*$. Тогда из свойств решения задачи (3.120) следует, что

$$\lim_{t \rightarrow \infty} \frac{\partial x}{\partial t} = 0,$$

и потому x^* является искомым решением для $Ax = b$.

При поиске значений $x(t)$, установившихся в пределе $t \rightarrow \infty$, значения $x(t)$ для конечных t не слишком интересны, так что для решения системы дифференциальных уравнений (3.120) можно применить простейший явный *метод Эйлера* (метод ломаных) с постоянным временным шагом τ , в котором производная заменяется на разделённую разность вперёд. Обозначая $x^{(k)} := x(t_k)$, $t_k = t_0 + \tau k$, $k = 0, 1, 2, \dots$, получим вместо (3.120)

$$\frac{x^{(k+1)} - x^{(k)}}{\tau} + Ax^{(k)} = b, \quad (3.121)$$

или

$$x^{(k+1)} = x^{(k)} - \tau(Ax^{(k)} - b), \quad k = 0, 1, 2, \dots,$$

то есть известный нам метод простой итерации (3.98) для решения системы уравнений $Ax = b$. При переменном шаге по времени, когда $\tau = \tau_k$, $k = 0, 1, 2, \dots$, получающийся метод Эйлера

$$\frac{x^{(k+1)} - x^{(k)}}{\tau_k} + Ax^{(k)} = b$$

эквивалентен

$$x^{(k+1)} = x^{(k)} - \tau_k(Ax^{(k)} - b), \quad k = 0, 1, 2, \dots,$$

т. е. простейшему нестационарному итерационному методу Ричардсона (3.110).

Представление метода простой итерации в виде (3.121), как метода решения системы дифференциальных уравнений, даёт возможность понять суть ограничений на параметр τ . Это не что иное, как ограничение на величину шага по времени, вызванное требованием устойчивости метода. Если шаг по времени мал, то до установления решения задачи (3.120) нам нужно сделать большое количество таких мелких шагов, что даёт ещё одно объяснение невысокой вычислительной эффективности метода простой итерации.

Более быструю сходимость к решению можно достичь, взяв шаг по времени большим, но для этого нужно преодолеть ограничение на устойчивость метода. Реализация этой идеи действительно приводит к более эффективным численным методам решения некоторых специальных систем линейных уравнений $Ax = b$, встречающихся при дискретизации дифференциальных уравнений с частными производными. Таковы *методы переменных направлений, методы расщепления и методы дробных шагов*, идейно близкие друг другу (см. [85]).

Очевидно, что вместо (3.120) можно рассмотреть задачу более общего вида

$$B \frac{\partial x}{\partial t} + Ax(t) = b, \quad (3.122)$$

где B — некоторая неособенная матрица. Смысл её введения станет более понятен, если переписать (3.122) в равносильном виде

$$\frac{\partial x}{\partial t} + B^{-1}Ax(t) = B^{-1}b.$$

Тогда в пределе, при занулении $\partial x/\partial t$, имеем

$$B^{-1}Ax = B^{-1}b.$$

откуда видно, что матрица B выполняет роль, аналогичную роли преобуславливающей матрицы для системы $Ax = b$.

Отметим в заключение темы, что для решения систем линейных алгебраических уравнений, возникающих при дискретизации уравнений в частных производных эллиптического типа, предельно эффективными являются так называемые многосеточные методы, предложенные Р.П. Федоренко в начале 60-х годов XX века.

3.12 Теория А.А. Самарского

Мы уже отмечали, что системы линейных алгебраических уравнений, которые необходимо решать на практике, часто бывают заданы неявно, в операторном виде. При этом мы не можем оперировать итерационными формулами вида (3.91) с явно заданным оператором T_k (наподобие (3.92)). Для подобных случаев А.А. Самарским была предложена специальная каноническая форма одношагового линейного итерационного процесса, предназначенного для решения систем уравнений $Ax = b$:

$$B_k \frac{x^{(k+1)} - x^{(k)}}{\tau_k} + Ax^{(k)} = b, \quad k = 0, 1, 2, \dots, \quad (3.123)$$

где B_k , τ_k — некоторые последовательности матриц и скалярных параметров соответственно, причём $\tau_k > 0$. Мы будем называть её *канонической формой Самарского*. Если $x^{(k)}$ сходится к пределу, то при некоторых необременительных условиях на B_k и τ_k этот предел является решением системы линейных алгебраических уравнений $Ax = b$.

Различные последовательности матриц B_k и итерационных параметров τ_k задают различные итерационные методы. Выбирая начальное значение $x^{(0)}$, находим затем из (3.123) последовательные приближения как решения уравнений

$$B_k x^{(k+1)} = (B_k - \tau_k I) x^{(k)} + \tau_k b, \quad k = 0, 1, 2, \dots$$

Ясно, что для однозначной разрешимости этой системы уравнений все матрицы B_k должны быть неособенными. Итерационный метод в фор-

ме (3.123) естественно назвать *явным*, если $B_k = I$ — единичная матрица и выписанная выше система сводится к явной формуле для нахождения следующего итерационного приближения $x^{(k+1)}$. Иначе, если $B_k \neq I$, итерации (3.123) называются *неявными*. Неявные итерационные методы имеют смысл применять лишь в том случае, когда решение системы уравнений относительно $x^{(k+1)}$ существенно легче, чем решение исходной системы.

Выпишем представление в форме Самарского для рассмотренных ранее итерационных процессов. Метод простой итерации из §3.9г принимает вид

$$\frac{x^{(k+1)} - x^{(k)}}{\tau} + Ax^{(k)} = b, \quad k = 0, 1, 2, \dots, \quad (3.124)$$

где $\tau = \tau_k = \text{const}$ — постоянный параметр, имеющий тот же смысл, что и в рассмотренных §3.9г. Переменный параметр τ_k в (3.124) приводит к нестационарному методу Ричардсона (3.110) (см. §3.10а). Если D и \tilde{L} — диагональная и строго нижняя треугольная части матрицы A соответственно (см. §3.9д), то методы Якоби и Гаусса-Зейделя можно записать в виде

$$D \frac{x^{(k+1)} - x^{(k)}}{1} + Ax^{(k)} = b,$$

и

$$(D + \tilde{L}) \frac{x^{(k+1)} - x^{(k)}}{1} + Ax^{(k)} = b.$$

Наконец, итерационный метод релаксации с релаксационным параметром ω (см. §3.9ж) в тех же обозначениях имеет форму Самарского

$$(D + \omega\tilde{L}) \frac{x^{(k+1)} - x^{(k)}}{\omega} + Ax^{(k)} = b, \quad k = 0, 1, 2, \dots$$

При исследовании сходимости итераций в форме Самарского удобно пользоваться матричными неравенствами, связанными со знакоопределённостью матриц. Условимся для вещественной $n \times n$ -матрицы G писать $G \triangleright 0$, если $\langle Gx, x \rangle > 0$ для всех ненулевых n -векторов x , т. е. если матрица G положительно определена. Из этого неравенства следует также существование такой константы $\mu > 0$, что $\langle Gx, x \rangle > \mu \langle x, x \rangle$. Неравенство $G \triangleright H$ будем понимать как $\langle Gx, x \rangle > \langle Hx, x \rangle$ для всех x , что равносильно также $G - H \triangleright 0$.

Достаточное условие сходимости итерационного процесса в форме Самарского (3.123) даёт

Теорема 3.12.1 (теорема Самарского) *Если A — симметричная положительно определённая матрица, $\tau > 0$ и $B \triangleright \frac{1}{2} \tau A$, то стационарный итерационный процесс*

$$B \frac{x^{(k+1)} - x^{(k)}}{\tau} + Ax^{(k)} = b, \quad k = 0, 1, 2, \dots,$$

сходится к решению системы уравнений $Ax = b$ из любого начального приближения.

Доказательство. Пусть x^* — решение системы уравнений $Ax = b$, так что

$$B \frac{x^* - x^*}{\tau} + Ax^* = b.$$

Если обозначить через $z^{(k)} = x^{(k)} - x^*$ — погрешность k -го приближения, то она удовлетворяет однородному соотношению

$$B \frac{z^{(k+1)} - z^{(k)}}{\tau} + Az^{(k)} = 0, \quad k = 0, 1, 2, \dots \quad (3.125)$$

Исследуем поведение энергетической нормы погрешности. Покажем сначала, что в условиях теоремы числовая последовательность $\|z^{(n)}\|_A = \langle Az^{(n)}, z^{(n)} \rangle$ является невозрастающей.

Из соотношения (3.125) следует

$$z^{(k+1)} = (I - \tau B^{-1}A) z^{(k)}, \quad (3.126)$$

и

$$Az^{(k+1)} = (A - \tau AB^{-1}A) z^{(k)}.$$

Таким образом,

$$\begin{aligned} \langle Az^{(k+1)}, z^{(k+1)} \rangle &= \langle Az^{(k)}, z^{(k)} \rangle - \tau \langle AB^{-1}Az^{(k)}, z^{(k)} \rangle \\ &\quad - \tau \langle Az^{(k)}, B^{-1}Az^{(k)} \rangle + \tau^2 \langle AB^{-1}Az^{(k)}, AB^{-1}Az^{(k)} \rangle. \end{aligned}$$

Коль скоро матрица A симметрична,

$$\langle AB^{-1}Az^{(k)}, z^{(k)} \rangle = \langle Az^{(k)}, B^{-1}Az^{(k)} \rangle,$$

и потому

$$\begin{aligned} \langle Az^{(k+1)}, z^{(k+1)} \rangle &= \\ &= \langle Az^{(k)}, z^{(k)} \rangle - 2\tau \langle (B - \frac{1}{2} \tau A) B^{-1}Az^{(k)}, B^{-1}Az^{(k)} \rangle. \quad (3.127) \end{aligned}$$

Учитывая неравенство $B \triangleright \frac{1}{2} \tau A$, можем заключить, что вычитаемое в правой части полученного равенства всегда неотрицательно. По этой причине

$$\|z^{(k+1)}\|_A \leq \|z^{(k)}\|_A,$$

так что последовательность $\|z^{(k)}\|_A$ монотонно не возрастает и ограничена снизу нулём. В силу теоремы Вейерштрасса она имеет предел при $k \rightarrow \infty$.

Неравенство $B \triangleright \frac{1}{2} \tau A$, т. е. положительная определённость матрицы $(B - \frac{1}{2} \tau A)$, означает существование такого $\eta > 0$, что для любых $y \in \mathbb{R}^n$

$$\langle (B - \frac{1}{2} \tau A) y, y \rangle \geq \eta \langle y, y \rangle = \eta \|y\|_2^2.$$

Окончательно получаем из (3.127)

$$\|z^{(k+1)}\|_A^2 - \|z^{(k)}\|_A^2 + 2\eta\tau \|B^{-1}Az^{(k)}\|_2^2 \leq 0$$

для всех $k = 0, 1, 2, \dots$. Переходя в этом неравенстве к пределу по $k \rightarrow \infty$, заключаем, что тогда $\|B^{-1}Az^{(k)}\|_2 \rightarrow 0$. При неособенной матрице $B^{-1}A$ это возможно лишь при $z^{(k)} \rightarrow 0$. Итак, вне зависимости от выбора начального приближения итерационный процесс в самом деле сходится. ■

Отметим, что из теоремы Самарского следует теорема Островского-Райха (Теорема 3.9.3) о сходимости метода релаксации для СЛАУ с симметричными положительно определёнными матрицами, а также, как её частный случай, Теорема 3.9.2 о сходимости метода Гаусса-Зейделя. В самом деле, пусть $A = \tilde{L} + D + \tilde{U}$ в обозначениях §3.9д, т. е. \tilde{L} и \tilde{U} — строго нижняя и строго верхняя треугольные части матрицы A , а D — её диагональная часть. Если A симметрична, то $\tilde{L} = \tilde{U}^\top$, и поэтому

$$\langle Ax, x \rangle = \langle \tilde{L}x, x \rangle + \langle Dx, x \rangle + \langle \tilde{U}x, x \rangle = \langle Dx, x \rangle + 2\langle \tilde{L}x, x \rangle.$$

Тогда

$$\begin{aligned} \langle Bx, x \rangle - \frac{1}{2} \omega \langle Ax, x \rangle &= \langle (D + \omega \tilde{L})x, x \rangle - \frac{1}{2} \omega (\langle Dx, x \rangle + 2\langle \tilde{L}x, x \rangle) \\ &= (1 - \frac{1}{2} \omega) \langle Dx, x \rangle > 0 \end{aligned}$$

при $0 < \omega < 2$.

Дальнейшие результаты в этом направлении читатель может увидеть, к примеру, в [37, 80].

3.13 Вычисление определителей матриц и обратных матриц

Предположим, что для матрицы A выполняется LU-разложение. Как отмечалось, выполняемые в представленной нами версии метода Гаусса преобразования — линейное комбинирование строк — не изменяют величины определителя матрицы. Следовательно, $\det A$ равен определителю получающейся в итоге верхней треугольной матрицы U , т. е. $\det A$ есть произведение диагональных элементов U .

Другая возможная трактовка этого результата состоит в том, что если $A = LU$ — треугольное разложение матрицы A , то, как известно из линейной алгебры,

$$\det A = \det L \cdot \det U.$$

Определитель нижней треугольной матрицы L равен 1, коль скоро на её диагонали стоят все единицы. Следовательно, как и ранее, $\det A = \det U$, а точнее — произведению всех диагональных элементов в верхней треугольной матрице U .

Совершенно аналогичные выводы можно сделать и при использовании других матричных разложений. Например, если нам удалось получить $A = QR$ — разложение исходной матрицы в произведение ортогональной и правой треугольной, то, коль скоро $\det Q = 1$, искомым определителем $\det A = \det R$ и вычисляется по R как произведение её диагональных элементов.

Рассмотрим теперь вычисление матрицы, обратной к данной матрице. Отметим, прежде всего, что в современных вычислительных технологиях это приходится делать не слишком часто. Один из примеров, когда подобное вычисление необходимо по существу, — нахождение дифференциала операции обращения матрицы $A \mapsto A^{-1}$, равного

$$d(A^{-1}) = -A^{-1}(dA)A^{-1}.$$

(см., к примеру, [14]). Тогда коэффициенты чувствительности решения системы уравнений $Ax = b$ по отношению к элементам матрицы и правой части (т. е. производные решения по коэффициентам и правым частям системы, см. §1.3) даются формулами

$$\frac{\partial x_\nu}{\partial a_{ij}} = -z_{\nu i} x_j, \quad \frac{\partial x_\nu}{\partial b_i} = z_{\nu i},$$

$\nu = 1, 2, \dots, n$, где $Z = (z_{ij}) = A^{-1}$ — обратная к матрице A .

Гораздо чаще встречается необходимость вычисления произведения обратной матрицы A^{-1} на какой-то вектор b , и это произведение всегда следует находить как решение системы уравнений $Ax = b$ какими-либо из методов для решения СЛАУ. Такой способ заведомо лучше, чем вычисление $A^{-1}b$ через нахождение обратной A^{-1} , как по точности, так и по трудоёмкости.

Матрица A^{-1} , обратная к данной матрице A , является решением матричного уравнения

$$AX = I.$$

Но это уравнение распадается на n уравнений относительно векторных неизвестных, соответствующих отдельным столбцам неизвестной матрицы X , и потому мы можем решать получающиеся уравнения по-рознь.

Из сказанного следует способ нахождения обратной матрицы: нужно решить n штук систем линейных уравнений

$$Ax = e^{(j)}, \quad j = 1, 2, \dots, n, \quad (3.128)$$

где $e^{(j)}$ — j -ый столбец единичной матрицы I . Это можно сделать, к примеру, любым из рассмотренных нами выше методов, причём прямые методы здесь особенно удобны в своей матричной трактовке. В самом деле, сначала мы можем выполнить один раз LU-разложение (или QR-разложение) исходной матрицы A , а затем хранить его и использовать посредством схемы (3.59) (или (3.77)) для различных правых частей уравнений (3.128). Если матрица A — симметричная положительно определённая, то очень удобным может быть разложение Холецкого и последующее решение систем уравнений (3.128) с помощью представления (3.67).

В прямых методах решения СЛАУ прямой ход, т. е. приведение исходной системы к треугольному виду, является наиболее трудоёмкой частью всего алгоритма, которая требует обычно $O(n^3)$ арифметических операций. Обратный ход (обратная подстановка) — существенно более лёгкая часть алгоритма, требующая всего $O(n^2)$ операций. По этой причине изложенный выше рецепт однократного LU-разложения матрицы (или других разложений) позволяет сохранить общую трудоёмкость $O(n^3)$ для алгоритма вычисления обратной матрицы.

Другой подход к обращению матриц — конструирование чисто матричных процедур, не опирающихся на методы решения систем линейных уравнений с векторными неизвестными. Известен итерационный

метод Шульца для обращения матриц: задавшись специальным начальным приближением $X^{(0)}$, выполняют итерации

$$X^{(k+1)} \leftarrow X^{(k)} (2I - AX^{(k)}), \quad k = 0, 1, 2, \dots \quad (3.129)$$

Метод Шульца — это не что иное как метод Ньютона для решения системы уравнений, применённый к $X^{-1} - A = 0$ (см. §4.5б).²⁵ Его можно также рассматривать как матричную версию известной процедуры для вычисления обратной величины (см. [12], глава 3).

Предложение 3.13.1 *Метод Шульца сходится тогда и только тогда, когда его начальное приближение $X^{(0)}$ удовлетворяет условию $\rho(I - AX^{(0)}) < 1$.*

Доказательство. Расчётную формулу метода Шульца можно переписать в виде

$$X^{(k+1)} = 2X^{(k)} - X^{(k)}AX^{(k)}.$$

Умножим обе части этого равенства слева на $(-A)$ и добавим к ним по единичной матрице I , получим

$$I - AX^{(k+1)} = I - 2AX^{(k)} + AX^{(k)}AX^{(k)},$$

что равносильно

$$I - AX^{(k+1)} = (I - AX^{(k)})^2, \quad k = 0, 1, 2, \dots$$

Отсюда, в частности, следует, что

$$I - AX^{(k)} = (I - AX^{(0)})^{2^k}, \quad k = 0, 1, 2, \dots$$

Если $X^{(k)} \rightarrow A^{-1}$ при $k \rightarrow \infty$, то $(I - AX^{(0)})^{2^k} \rightarrow 0$ — последовательность степеней матрицы сходится к нулю. Тогда необходимо $\rho(I - AX^{(0)}) < 1$ в силу Предложения 3.3.10.

И наоборот, если $\rho(I - AX^{(0)}) < 1$, то $(I - AX^{(0)})^{2^k} \rightarrow 0$ при $k \rightarrow \infty$, и потому должна иметь место сходимость $X^{(k)} \rightarrow A^{-1}$. ■

Из доказательства предложения следует, что метод Шульца имеет квадратичную сходимость.

²⁵Иногда этот метод называют также *методом Хотеллинга*, так как одновременно с Г. Шульцем [94] его рассматривал американский экономист и статистик Г. Хотеллинг [89]. Кроме того, встречается (хотя и крайне редко) также название *метод Бодевига*.

3.14 Оценка погрешности приближённого решения

В этом параграфе мы рассмотрим практически важный вопрос об оценке погрешности приближённого решения систем линейных алгебраических уравнений. Первый способ носит общий характер и может применяться в любых ситуациях, в частности, не обязательно в связи с итерационными методами.

Пусть \tilde{x} — приближённое решение системы уравнений $Ax = b$, тогда как x^* — её точное решение. Тогда, принимая во внимание, что $I = A^{-1}A$ и $Ax^* = b$,

$$\begin{aligned} \|\tilde{x} - x^*\| &= \|A^{-1}A\tilde{x} - A^{-1}Ax^*\| \\ &= \|A^{-1}(A\tilde{x} - Ax^*)\| \\ &\leq \|A^{-1}\| \|A\tilde{x} - b\|, \end{aligned} \quad (3.130)$$

где матричная и векторная нормы, естественно, должны быть согласованы. Величина $(A\tilde{x} - b)$ — это невязка приближённого решения \tilde{x} , которую мы обычно можем вычислять непосредственно по \tilde{x} . Как следствие, погрешность решения можно узнать, найдя каким-либо образом или оценив сверху норму обратной матрицы $\|A^{-1}\|$.

Иногда из практики можно получать какую-то информацию о значении $\|A^{-1}\|$. Например, если A — симметричная положительно определённая матрица и известна нижняя граница её спектра $\mu > 0$, то $\|A^{-1}\|_2 \leq 1/\mu$. Напомним, что аналогичную информацию о спектре матрицы СЛАУ мы использовали при оптимизации скалярного предобуславливателя в §3.9г. Такова ситуация с численным решением некоторых популярных уравнений математической физики (уравнением Лапласа и его обобщениями, к примеру), для которых дискретные аналоги соответствующих дифференциальных операторов хорошо изучены и известны оценки их собственных значений.

В общем случае быстрое нахождение $\|A^{-1}\|$ или хотя бы разумных оценок в какой-то норме для $\|A^{-1}\|$ сверху, более быстрое, чем решение исходной СЛАУ, является нетривиальным делом. Краткий обзор существующих численных процедур для этой цели («оценщиков» обратной матрицы) и дальнейшие ссылки на литературу можно найти в [13].

Для конкретных численных методов оценка погрешности приближённого решения иногда может быть выведена из свойств этих ме-

тодов. Например, в стационарных одношаговых итерационных методах последовательность погрешностей приближений своими свойствами очень близка к геометрической прогрессии, и этим обстоятельством можно с успехом воспользоваться.

Пусть задан сходящийся стационарный одношаговый итерационный метод

$$x^{(k+1)} \leftarrow Cx^{(k)} + d, \quad k = 0, 1, 2, \dots,$$

в котором $\|C\| < 1$ для некоторой матричной нормы. Ясно, что ввиду результатов §3.9б о связи спектрального радиуса и матричных норм последнее допущение не ограничивает общности нашего рассмотрения. Как оценить отклонение по норме очередного приближения $x^{(k)}$ от предела $x^* := \lim_{k \rightarrow \infty} x^{(k)}$, не зная самого этого предела и наблюдая лишь за итерационной последовательностью $x^{(0)}, x^{(1)}, \dots, x^{(k)}, \dots$?

Как и прежде, имеем

$$\begin{aligned} x^{(k)} &= Cx^{(k-1)} + d, \\ x^* &= Cx^* + d. \end{aligned}$$

Вычитание второго равенства из первого даёт

$$x^{(k)} - x^* = C(x^{(k-1)} - x^*). \quad (3.131)$$

Перенесём $x^{(k)}$ в правую часть этого соотношения, а затем добавим к обеим частям по $x^{(k-1)}$:

$$x^{(k-1)} - x^* = x^{(k-1)} - x^{(k)} + C(x^{(k-1)} - x^*).$$

Возьмём теперь от обеих частей полученного равенства векторную норму, которая согласована с используемой матричной нормой для C . Применяя затем неравенство треугольника, приходим к оценке

$$\|x^{(k-1)} - x^*\| \leq \|x^{(k)} - x^{(k-1)}\| + \|C\| \cdot \|x^{(k-1)} - x^*\|,$$

Перенесение в левую часть второго слагаемого из правой части и последующее деление обеих частей неравенства на положительную величину $(1 - \|C\|)$ даёт

$$\|x^{(k-1)} - x^*\| \leq \frac{1}{1 - \|C\|} \|x^{(k)} - x^{(k-1)}\|. \quad (3.132)$$

С другой стороны, вспомним, что из (3.131) следует

$$\|x^{(k)} - x^*\| \leq \|C\| \cdot \|x^{(k-1)} - x^*\|.$$

Подставляя сюда вместо $\|x^{(k-1)} - x^*\|$ оценку сверху (3.132), получаем окончательно

$$\|x^{(k)} - x^*\| \leq \frac{\|C\|}{1 - \|C\|} \|x^{(k)} - x^{(k-1)}\|. \quad (3.133)$$

Выведенная оценка может быть использована на практике как для оценки погрешности какого-то приближения из итерационной последовательности, так и для определения момента окончания итераций, т. е. того, достигнута ли желаемая точность приближения к решению или нет.

Пример 3.14.1 Рассмотрим систему линейных алгебраических уравнений

$$\begin{pmatrix} 2 & 1 \\ 3 & 4 \end{pmatrix} x = \begin{pmatrix} 0 \\ 5 \end{pmatrix},$$

точное решение которой равно $(-1, 2)^\top$. Пусть для решения этой системы организован итерационный метод Гаусса-Зейделя с начальным приближением $x^{(0)} = (0, 0)^\top$. Через сколько итераций компоненты очередного приближения к решению станут отличаться от точного решения не более, чем на 10^{-3} ?

Исследуемый нами вопрос требует чебышёвской нормы $\|\cdot\|_\infty$ для измерения отклонения векторов друг от друга, и соответствующая подчинённая матричная норма задаётся выражением из Предложения 3.3.6. Матрица оператора перехода итерационного метода Гаусса-Зейделя согласно (3.103) есть

$$-\begin{pmatrix} 2 & 0 \\ 3 & 4 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & -0.5 \\ 0 & 0.375 \end{pmatrix},$$

так что её ∞ -норма равна 0.5. Следовательно, в оценке (3.133) имеем

$$\frac{\|C\|}{1 - \|C\|} = \frac{0.5}{1 - 0.5} = 1,$$

и потому должно быть справедливым неравенство

$$\|x^{(k)} - x^*\|_\infty \leq \|x^{(k)} - x^{(k-1)}\|_\infty. \quad (3.134)$$

Оно показывает, что компоненты очередного приближения отличаются от компонент точного решения не более, чем компоненты приближений друг от друга.

Запустив итерации Гаусса-Зейделя, мы можем видеть, что

$$\begin{aligned}x^{(0)} &= (0, 0)^\top, \\x^{(1)} &= (0, 1.25)^\top, \\x^{(2)} &= (-0.625, 1.71875)^\top, \\&\dots \quad \dots \\x^{(8)} &= (-0.998957, 1.999218)^\top, \\x^{(9)} &= (-0.999609, 1.999707)^\top,\end{aligned}$$

т. е. 9-я итерация отличается от предыдущей 8-й меньше чем на 10^{-3} , и потому согласно оценке (3.134) на этой итерации мы получаем требуемую погрешность. То, что она действительно такова, можно убедиться из сравнения $x^{(9)}$ с известным нам точным решением $(-1, 2)^\top$. ■

Как хорошо видно из примера, практическая реализация методики оценки погрешности итерационного решения может столкнуться с двумя трудностями. Во-первых, непростым является определение матрицы C (которая может и не задаваться в явном виде). Во-вторых, выбор нормы $\|\cdot\|$, в которой $\|C\| < 1$, также может быть неочевидным. Теоретически такая норма должна существовать, если итерационный процесс сходится из любого начального приближения, но её конкретный выбор в общем случае прост.

3.15 Линейная задача о наименьших квадратах

Для заданных $m \times n$ -матрицы A и m -вектора b *линейной задачей о наименьших квадратах* называют задачу отыскания такого вектора x , который доставляет минимум квадратичной форме $\langle Ax - b, Ax - b \rangle$, или, что равносильно, квадрату евклидовой нормы невязки $\|Ax - b\|_2^2$. Ясно, что для матриц A полного ранга в случае $m \leq n$, когда число строк матрицы не превосходит числа столбцов, искомым минимумом, как правило, равен нулю. Для квадратной матрицы A линейная задача о наименьших квадратах, фактически, равносильна решению системы

линейных алгебраических уравнений $Ax = b$ и несёт особую специфику лишь когда A имеет неполный ранг, т. е. особенна. Теоретически и практически наиболее важный случай линейной задачи о наименьших квадратах соответствует $m > n$. Он находит многочисленные и разнообразные применения при обработке данных

Коль скоро

$$\begin{aligned}\langle Ax - b, Ax - b \rangle &= \langle Ax, Ax \rangle - \langle b, Ax \rangle - \langle Ax, b \rangle + \langle b, b \rangle \\ &= \langle Ax, Ax \rangle - 2\langle Ax, b \rangle + \langle b, b \rangle,\end{aligned}$$

то

$$\frac{\partial}{\partial x_i} \langle Ax - b, Ax - b \rangle = \frac{\partial}{\partial x_i} (\langle Ax, Ax \rangle - 2\langle Ax, b \rangle)$$

Система линейных алгебраических уравнений

$$A^T A x = A^T b \quad (3.135)$$

называется *нормальной системой уравнений* для линейной задачи о наименьших квадратах с матрицей A и вектором b .²⁶ Её решение и доставляет искомым минимум выражению $\|Ax - b\|_2^2$

3.16 Проблема собственных значений

3.16а Обсуждение постановки задачи

Ненулевой вектор v называется *собственным вектором* квадратной матрицы A , если в результате умножения на эту матрицу он переходит в коллинеарный себе, т. е. отличающийся от исходного только некоторым скалярным множителем:

$$Av = \lambda v. \quad (3.136)$$

Сам скаляр λ , который является коэффициентом пропорциональности исходного вектора и его образа при действии матрицы, называют *собственным значением* матрицы. *Проблемой собственных значений* называют задачу определения собственных значений и собственных векторов матриц: для заданной $n \times n$ -матрицы A найти числа λ и n -векторы $v \neq 0$, удовлетворяющие условию (3.136).

²⁶Переход от исходной системы уравнений $Ax = b$ к нормальной системе (3.135) иногда называют *первой трансформацией Гаусса*.

Система уравнений (3.136) кажется недоопределённой, так как содержит $n + 1$ неизвестных, которые нужно найти из n уравнений. Но на самом деле можно замкнуть её, к примеру, каким-нибудь условием нормировки собственных векторов ($\|v\| = 1$ в какой-то норме) или требованием, чтобы какая-нибудь компонента v принимала бы заданное значение. Последнее условие иногда даже более предпочтительно ввиду своей линейности.

Даже если рассматриваемая матрица A имеет все элементы вещественными, могут не существовать вещественные λ и v , удовлетворяющие соотношению (3.136). По этой причине целесообразно рассматривать задачу определения собственных чисел λ и собственных векторов v в поле комплексных чисел \mathbb{C} , которое алгебраически замкнуто.²⁷ Из курса линейной алгебры читателю должно быть известно, что задача нахождения собственных значений матрицы A эквивалентна задаче нахождения корней уравнения

$$\det(A - \lambda I) = 0,$$

называемого *характеристическим* (или *вековым*) уравнением для матрицы A .

Иногда при упоминании этой задачи подчёркивают — «алгебраическая проблема собственных значений», чтобы уточнить, что речь идёт о матрицах конечных размеров, конечномерной ситуации и т. п. в отличие, скажем, от задачи нахождения собственных значений операторов в бесконечномерных пространствах функций. Слово «проблема» также уместно в этом контексте, поскольку рассматриваемая задача сложна и имеет много аспектов.

Различают *полную проблему* собственных значений и *частичную проблему* собственных значений. В полной проблеме требуется нахождение всех собственных чисел и собственных векторов. Частичная проблема собственных значений — это задача нахождения некоторых собственных чисел матрицы и/или некоторых собственных векторов. К примеру, наибольшего по модулю собственного значения, или нескольких наибольших по модулю собственных значений и соответствующих им собственных векторов.

Собственные значения матриц нужно знать во многих приложениях. Например, задача определения частот собственных колебаний ме-

²⁷Напомним, что *алгебраически замкнутым* называется поле, в котором в котором всякий многочлен ненулевой степени с коэффициентами из этого поля имеет хотя бы один корень.

ханических систем (весьма важная, к примеру, при проектировании различных конструкций) сводится к нахождению собственных значений так называемых матриц жёсткости этих систем. Особую важность собственным значениям придаёт то обстоятельство, что соответствующие им частоты собственных колебаний являются непосредственно наблюдаемыми из опыта физическими величинами. Это тон звучания тронутой гитарной струны и т. п.

Пример 3.16.1 Линейные динамические системы с дискретным временем вида

$$x^{(k+1)} = Ax^{(k)} + b^{(k)}, \quad k = 0, 1, 2, \dots, \quad (3.137)$$

служат моделями разнообразных процессов окружающего нас мира, от биологии до экономики.

Общее решение такой системы есть сумма частного решения исходной системы (3.137) и общего решения однородной системы $x^{(k+1)} = Ax^{(k)}$ без свободного члена. Если искать нетривиальные решения однородной системы в виде $x^{(k)} = \lambda^k h$, где λ — ненулевой скаляр и h — n -вектор, то нетрудно убедиться, что λ должно быть собственным значением A , а h — собственным вектором матрицы A . ■

Ясно, что собственные векторы матрицы определяются неоднозначно, с точностью до скалярного множителя. В связи с этим часто говорят о нахождении одномерных *инвариантных подпространств* матрицы. Инвариантные подпространства матрицы могут иметь и большую размерность, и в любом случае их знание доставляет важную информацию о рассматриваемом линейном операторе, позволяя упростить его представление. Пусть, например, \mathcal{S} — это p -мерное инвариантное подпространство матрицы A , так что $Ax \in \mathcal{S}$ для любого $x \in \mathcal{S}$, и базисом \mathcal{S} являются векторы v_1, v_2, \dots, v_p . Беря базис всего пространства \mathbb{R}^n так, чтобы его последними векторами были v_1, v_2, \dots, v_p (это, очевидно, можно сделать всегда), получим в нём блочно-треугольное представление рассматриваемого линейного оператора:

$$\begin{pmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{pmatrix}$$

с $p \times p$ -блоком A_{22} . В последние десятилетия задача определения для матрицы тех или иных инвариантных подпространств, не обязательно одномерных, также включается в «проблему собственных значений».

Помимо необходимости выхода в общем случае в комплексную плоскость \mathbb{C} , даже для вещественных матриц, ещё одной особенностью проблемы собственных значений, осложняющей её решение является нелинейный характер задачи, несмотря на традицию отнесения её к «вычислительной линейной алгебре». Это обстоятельство нетрудно осознать из рассмотрения основного соотношения (3.136)

$$Av = \lambda v,$$

которое является системой уравнений относительно λ и v , причём в его правой части суммарная степень неизвестных переменных равна *двум*: $2 = (1 \text{ при } \lambda) + (1 \text{ при } v)$.

Если собственное значение $\tilde{\lambda}$ матрицы A уже найдено, то, как известно, определение соответствующих собственных векторов сводится к решению системы линейных алгебраических уравнений

$$(A - \tilde{\lambda}I)x = 0$$

с особенной матрицей. Но на практике часто предпочитают пользоваться для нахождения собственных векторов специализированными вычислительными процедурами. Многие из них позволяют вычислять собственные векторы одновременно с собственными значениями матриц.

В заключение нашего обсуждения коснёмся алгоритмического аспекта проблемы собственных значений. Напомним известную в алгебре теорему Абея-Руффини: для алгебраических полиномов степени выше 4 не существует прямых методов нахождения корней. Как следствие, мы не вправе ожидать существования прямых методов решения проблемы собственных значений для произвольных матриц размера более 4×4 , и потому рассматриваемые ниже методы — существенно итерационные.

3.166 Обусловленность проблемы собственных значений

Спектр матрицы, как множество точек комплексной плоскости \mathbb{C} , непрерывно зависит от элементов матрицы. Соответствующий результат часто называют теоремой Островского (читатель может увидеть детальное изложение этой теории в книгах [19, 26, 34, 41, 50]). Но собственные векторы (инвариантные подпространства) матрицы могут из-

меняться в зависимости от матрицы разрывным образом даже в совершенно обычных ситуациях.

Пример 3.16.2 [50] Рассмотрим матрицу

$$A = \begin{pmatrix} 1 + \varepsilon & \delta \\ 0 & 1 \end{pmatrix}.$$

Её собственные значения суть числа 1 и $1 + \varepsilon$, и при $\varepsilon\delta \neq 0$ соответствующими нормированными собственными векторами являются

$$\frac{1}{\sqrt{\varepsilon^2 + \delta^2}} \begin{pmatrix} -\delta \\ \varepsilon \end{pmatrix} \quad \text{и} \quad \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Выбирая подходящим образом отношение ε/δ , можно придать первому собственному вектору любое направление, сколь бы малыми ни являлись значения ε и δ .

Если положить $\varepsilon = 0$, то

$$A = \begin{pmatrix} 1 & \delta \\ 0 & 1 \end{pmatrix}.$$

При $\delta \neq 0$ у матрицы A будет всего один собственный вектор, хотя при надлежащем δ её можно сделать сколь угодно близкой к единичной матрице, имеющей два линейно независимых собственных вектора. ■

При более пристальном изучении проблемы собственных значений выясняется, что, несмотря на непрерывную зависимость собственных значений от элементов матрицы, скорость их изменения может быть сколь угодно большой (даже для матриц фиксированного размера), если они соответствуют так называемым нелинейным элементарным делителям матрицы — жордановым клеткам размера 2 и более.

Пример 3.16.3 Собственные значения матрицы

$$A = \begin{pmatrix} \lambda & 1 \\ \varepsilon & \lambda \end{pmatrix}$$

— возмущённой жордановой 2×2 -клетки — равны $\lambda \pm \sqrt{\varepsilon}$, так что мгновенная скорость их изменения, равная $\sqrt{\varepsilon}/\varepsilon$, при $\varepsilon = 0$ бесконечна.

Это же явление имеет место и для произвольной жордановой клетки, размером более двух. ■

Всюду далее большую роль будут играть матрицы, которые преобразованием подобия можно привести к диагональному виду. Для их обозначения вводится

Определение 3.16.1 *Матрицы, подобные диагональным матрицам, будем называть матрицами простой структуры или диагонализуемыми матрицами.*²⁸

Собственные числа матриц простой структуры зависят от возмущений гораздо более «плавным образом», чем в общем случае.

Теорема 3.16.1 (теорема Бауэра-Файка [86]) *Если A — квадратная матрица простой структуры, $\lambda_i(A)$ — её собственные числа, V — матрица из собственных векторов A , а $\tilde{\lambda}$ — собственное число возмущённой матрицы $A + \Delta A$, то*

$$\min_i |\tilde{\lambda} - \lambda_i(A)| \leq \text{cond}_2(V) \|\Delta A\|_2. \quad (3.138)$$

Доказательство. Если $\tilde{\lambda}$ совпадает с каким-то из собственных значений исходной матрицы A , то левая часть доказываемого неравенства зануляется, и оно, очевидно, справедливо. Будем поэтому предполагать, что $\tilde{\lambda}$ не совпадает ни с одним из $\lambda_i(A)$, $i = 1, 2, \dots, n$. Следовательно, если, согласно условию теоремы

$$V^{-1}AV = D,$$

где $D = \text{diag} \{ \lambda_1, \lambda_2, \dots, \lambda_n \}$ — диагональная матрица с собственными числами матрицы A по диагонали, то матрица $D - \tilde{\lambda}I$ неособенна.

С другой стороны, матрица $A + \Delta A - \tilde{\lambda}I$ является особенной по построению, так что особенна и матрица $V^{-1}(A + \Delta A - \tilde{\lambda}I)V$. Но

$$\begin{aligned} V^{-1}(A + \Delta A - \tilde{\lambda}I)V &= (D - \tilde{\lambda}I) + V^{-1}(\Delta A)V \\ &= (D - \tilde{\lambda}I)(I + (D - \tilde{\lambda}I)^{-1}V^{-1}(\Delta A)V), \end{aligned}$$

²⁸Такие матрицы называют также *недефектными*.

откуда можно заключить о том, что особенной должна быть матрица $(I + (D - \tilde{\lambda}I)^{-1}V^{-1}(\Delta A)V)$. Как следствие, матрица

$$(D - \tilde{\lambda}I)^{-1}V^{-1}(\Delta A)V$$

имеет собственное значение -1 , и потому любая норма этой матрицы должна быть не меньшей 1. В частности, это верно для спектральной нормы:

$$\|(D - \tilde{\lambda}I)^{-1}V^{-1}(\Delta A)V\|_2 \geq 1.$$

Отсюда

$$\max_{1 \leq i \leq n} |(\lambda_i - \tilde{\lambda})^{-1}| \cdot \|V^{-1}\|_2 \|\Delta A\|_2 \|V\|_2 \geq 1.$$

Последнее неравенство равносильно

$$\min_{1 \leq i \leq n} |(\lambda_i - \tilde{\lambda})^{-1}| \leq \|V^{-1}\|_2 \|\Delta A\|_2 \|V\|_2,$$

или

$$\min_i |\tilde{\lambda} - \lambda_i(A)| \leq \text{cond}_2(V) \|\Delta A\|_2,$$

как и требовалось. ■

Теорема Бауэра-Файка показывает, что, каково бы ни было возмущение ΔA матрицы простой структуры A , для любого собственного значения $\tilde{\lambda}$ возмущённой матрицы $A + \Delta A$ найдётся собственное значение λ_i матрицы A , отличающееся от $\tilde{\lambda}$ не более чем на величину спектральной нормы возмущения $\|\Delta A\|_2$, умноженную на число обусловленности матрицы собственных векторов. Таким образом, число обусловленности матрицы из собственных векторов может служить мерой обусловленности проблемы собственных значений.

Практическую ценность теоремы Бауэра-Файка в целом и неравенства (3.138) в частности снижает то обстоятельство, что собственные векторы матрицы определены с точностью до скалярного множителя, и потому $\text{cond}_2(V)$ есть величина, заданная не вполне однозначно. Наилучшим выбором для $\text{cond}_2(V)$ в неравенстве был бы, очевидно, минимум чисел обусловленности матриц из собственных векторов, но его нахождение является в общем случае сложной задачей. Тем не менее, прикидочные оценки и качественные выводы на основе теоремы Бауэра-Файка делать можно.

Важнейший частный случай применения теоремы Бауэра-Файка относится к симметричным матрицам. Они имеют простую структуру и,

кроме того, собственные векторы симметричных матриц ортогональны друг другу. Как следствие, матрица собственных векторов V может быть взята ортогональной, с числом обусловленности 1. Получаем следующий результат: если $\lambda_i(A)$ — собственные числа симметричной матрицы A , а $\tilde{\lambda}$ — собственное число возмущённой матрицы $A + \Delta A$, то

$$\min_i |\tilde{\lambda} - \lambda_i(A)| \leq \|\Delta A\|_2.$$

Иными словами, при возмущении симметричных матриц их собственные числа изменяются на величину, не превосходящую спектральной нормы возмущения, т. е. гораздо более умеренно, чем для матриц общего вида.

Предложение 3.16.1 *Матрицы простой структуры образуют открытое всюду плотное подмножество во множестве всех квадратных матриц.*

Набросок доказательства таков: если для произвольного малого ε к канонической жордановой форме $n \times n$ -матрицы прибавить возмущающую матрицу вида $\text{diag} \{ \varepsilon, \varepsilon/2, \varepsilon/3, \dots, \varepsilon/n \}$, то получающаяся треугольная матрица будет иметь различные собственные числа, т. е. делается диагонализуемой. Требуемое возмущение исходной матрицы мы можем получить из $\text{diag} \{ \varepsilon, \varepsilon/2, \varepsilon/3, \dots, \varepsilon/n \}$ путём преобразования подобия, обратного по отношению к тому, которое переводит исходную матрицу к жордановой нормальной форме.

Как следствие, матрицы с нелинейными элементарными делителями, которые соответствуют жордановым клеткам размера 2 и более, составляют множество *первой бэровской категории* во множестве всех матриц. Подобные множества, называемые также *тощими*, являются в топологическом смысле наиболее разреженными и бедными множествами (см. [18, 48]). Но на долю таких матриц приходится главные трудности, с которыми сталкиваются при решении проблемы собственных значений. В этом отношении задача нахождения сингулярных чисел и сингулярных векторов является принципиально другой, так как симметричная матрица $A^T A$ (эрмитова матрица $A^* A$ в комплексном случае) всегда имеет простую структуру, т. е. диагонализуема.

3.16в Коэффициенты перекоса матрицы

Целью этого пункта является детальное исследование устойчивости решения проблемы собственных значений в упрощённой ситуации,

когда все собственные значения матрицы A различны. Именно в этом случае, как было отмечено в §3.16б, собственные векторы непрерывно зависят от элементов матрицы и, более того, существуют их конечные дифференциалы.

Пусть A — данная матрица и dA — дифференциал (главная линейная часть) её возмущения, так что $A + dA$ — это близкая к A возмущённая матрица. Как изменятся собственные значения и собственные векторы матрицы $A + dA$ в сравнении с собственными значениями и собственными векторами A ?

Имеем

$$\begin{aligned} Ax_i &= \lambda_i x_i, \\ (A + dA)(x_i + dx_i) &= (\lambda_i + d\lambda_i)(x_i + dx_i), \end{aligned}$$

где через λ_i обозначены собственные значения A , x_i — собственные векторы, $i = 1, 2, \dots, n$, причём последние образуют базис в \mathbb{R}^n , коль скоро по предположению A является матрицей простой структуры. Пренебрегая членами второго порядка малости, можем выписать равенство

$$(dA)x_i + A(dx_i) = \lambda_i(dx_i) + (d\lambda_i)x_i. \quad (3.139)$$

Пусть y_1, y_2, \dots, y_n — собственные векторы эрмитово-сопряжённой матрицы A^* , соответствующие её собственным значениям $\bar{\lambda}_1, \bar{\lambda}_2, \dots, \bar{\lambda}_n$. Умножая скалярно равенство (3.139) на y_j , получим

$$\langle (dA)x_i, y_j \rangle + \langle A(dx_i), y_j \rangle = \lambda_i \langle dx_i, y_j \rangle + (d\lambda_i) \langle x_i, y_j \rangle. \quad (3.140)$$

В частности, при $j = i$ имеем

$$\langle (dA)x_i, y_i \rangle + \langle A(dx_i), y_i \rangle = \lambda_i \langle dx_i, y_i \rangle + (d\lambda_i) \langle x_i, y_i \rangle,$$

где соседние со знаком равенства члены можно взаимно уничтожить: они оказываются одинаковыми, коль скоро

$$\langle A(dx_i), y_i \rangle = \langle dx_i, A^*y_i \rangle = \langle dx_i, \bar{\lambda}_i y_i \rangle = \lambda_i \langle dx_i, y_i \rangle.$$

Следовательно,

$$\langle (dA)x_i, y_i \rangle = (d\lambda_i) \langle x_i, y_i \rangle,$$

и потому

$$d\lambda_i = \frac{\langle (dA)x_i, y_i \rangle}{\langle x_i, y_i \rangle}.$$

Пусть теперь $j \neq i$. Тогда $\langle x_i, y_j \rangle = 0$ в силу биортогональности систем векторов $\{x_i\}$ и $\{y_j\}$ (Предложение 3.2.2), и потому

$$\langle A(dx_i), y_j \rangle = \langle dx_i, A^* y_j \rangle = \langle dx_i, \bar{\lambda}_j y_j \rangle = \lambda_j \langle dx_i, y_j \rangle.$$

Подставляя этот результат в (3.140), будем иметь

$$\langle (dA)x_i, y_j \rangle + \lambda_j \langle dx_i, y_j \rangle = \lambda_i \langle dx_i, y_j \rangle.$$

Поэтому

$$\langle dx_i, y_j \rangle = \frac{\langle (dA)x_i, y_j \rangle}{\lambda_i - \lambda_j}.$$

Разложим dx_i по базису из собственных векторов невозмущённой матрицы A :

$$dx_i = \sum_{j=1}^n \alpha_{ij} x_j.$$

Так как собственные векторы задаются с точностью до множителя, то в этом разложении коэффициенты α_{ii} содержательного смысла не имеют, и можно считать, что $\alpha_{ii} = 0$ (напомним, что мы, в действительности, ищем возмущение одномерного инвариантного подпространства матрицы). Для остальных коэффициентов имеем $\langle dx_i, y_j \rangle = \alpha_{ij} \langle x_j, y_j \rangle$, опять таки в силу Предложения 3.2.2. Следовательно, для $i \neq j$

$$\alpha_{ij} = \frac{\langle (dA)x_i, y_j \rangle}{(\lambda_i - \lambda_j) \langle x_j, y_j \rangle}.$$

Перейдём теперь к оцениванию возмущений собственных значений и собственных векторов. Из формулы для дифференциала $d\lambda_i$ и из неравенства Коши-Буняковского следует

$$|d\lambda_i| \leq \frac{\|dA\|_2 \|x\|_2 \|y\|_2}{\langle x_i, y_i \rangle} = \nu_i \|dA\|_2,$$

где посредством

$$\nu_i = \frac{\|x_i\|_2 \|y_i\|_2}{\langle x_i, y_i \rangle}, \quad i = 1, 2, \dots, n,$$

обозначены величины, называемые *коэффициентами перекоса* матрицы A , отвечающие собственным значениям λ_i , $i = 1, 2, \dots, n$.

Ясно, что $\nu_i \geq 1$, и можно интерпретировать коэффициенты перекоса как

$$\nu_i = \frac{1}{\cos \varphi_i},$$

где φ_i угол между собственными векторами x_i и y_i исходной и сопряжённой матриц. Коэффициенты перекоса характеризуют, таким образом, обусловленность проблемы собственных значений в смысле второго подхода §1.3.

Для симметричной (или, более общо, эрмитовой) матрицы коэффициенты перекоса равны 1. В самом деле, сопряжённая к ней задача на собственные значения совпадает с ней самой, и потому в наших обозначениях $x_i = y_i$, $i = 1, 2, \dots, n$. Следовательно, $\langle x_i, y_i \rangle = \langle x_i, x_i \rangle = \|x_i\|_2 \|y_i\|_2$, откуда и следует $\nu_i = 1$.

Это наименьшее возможное значение коэффициентов перекоса, так что численное нахождение собственных значений симметричных (эрмитовых в комплексном случае) матриц является наиболее устойчивым.

Что касается возмущений собственных векторов, то коэффициенты α_{ij} их разложения оцениваются сверху как

$$|\alpha_{ij}| \leq \frac{\|(dA)x_i\|_2 \|y_j\|_2}{|\lambda_i - \lambda_j| \cdot |\langle x_j, y_j \rangle|} \leq \frac{\|dA\|_2}{|\lambda_i - \lambda_j|} \nu_j,$$

и потому имеет место общая оценка

$$\|dx_i\|_2 \leq \|dA\|_2 \cdot \|x\|_2 \cdot \sum_{j \neq i} \frac{\nu_j}{|\lambda_i - \lambda_j|}. \quad (3.141)$$

Отметим значительную разницу в поведении возмущений собственных значений и собственных векторов матриц. Из оценки (3.141) следует, что на чувствительность отдельного собственного вектора влияют коэффициенты перекоса *всех* собственных значений матрицы, а не только того, которое отвечает этому вектору. Кроме того, в знаменателях слагаемых из правой части (3.141) присутствуют разности $\lambda_i - \lambda_j$, которые могут быть малыми при близких собственных значениях матрицы. Как следствие, собственные векторы при этом очень чувствительны к возмущениям в элементах матрицы. Это мы могли наблюдать в Примере 3.16.2. В частности, даже для симметричных (эрмитовых) матриц задача отыскания собственных векторов может оказаться плохообусловленной.

Рассмотрим l -ую компоненту векторного равенства (3.142):

$$\sum_{j=1}^n a_{lj}v_j = \lambda v_l,$$

что равносильно

$$\sum_{\substack{j=1 \\ j \neq l}}^n a_{lj}v_j = (\lambda - a_{ll})v_l.$$

По этой причине

$$\begin{aligned} |\lambda - a_{ll}| |v_l| &= \left| \sum_{j \neq l} a_{lj}v_j \right| \leq \sum_{j \neq l} |a_{lj}v_j| \\ &= \sum_{j \neq l} |a_{lj}| |v_j| \leq |v_l| \sum_{j \neq l} |a_{lj}|, \end{aligned}$$

коль скоро $|v_j| \leq |v_l|$. Наконец, поскольку $v \neq 0$, мы можем сократить обе части полученного неравенства на положительную величину $|v_l|$:

$$|\lambda - a_{ll}| \leq \sum_{j \neq l} |a_{lj}|.$$

Не зная собственного вектора v , мы не располагаем и номером l его наибольшей по модулю компоненты. Но можно действовать наверняка, рассмотрев дизъюнкцию соотношений выписанного выше вида для всех $l = 1, 2, \dots, n$, так как хотя бы для одного из них непременно справедливы выполненные нами рассуждения. Потому в целом, если λ — какое-либо собственное значение рассматриваемой матрицы A , должно выполняться хотя бы одно из неравенств

$$|\lambda - a_{ll}| \leq \sum_{j \neq l} |a_{lj}|, \quad l = 1, 2, \dots, n.$$

Каждое из этих соотношений на λ определяет на комплексной плоскости \mathbb{C} круг с центром в точке a_{ll} и радиусом, равным $\sum_{j \neq l} |a_{lj}|$. Как следствие, мы приходим к результату, который был установлен в 1931 году С.А. Гершгориним:

Теорема 3.16.2 (теорема Гершгорина) *Все собственные значения $\lambda(A)$ любой вещественной или комплексной $n \times n$ -матрицы $A = (a_{ij})$ расположены в объединении кругов комплексной плоскости с центрами a_{ii} и радиусами $\sum_{j \neq i} |a_{ij}|$, $i = 1, 2, \dots, n$, т. е.*

$$\lambda(A) \in \bigcup_{i=1}^n \left\{ z \in \mathbb{C} \mid |z - a_{ii}| \leq \sum_{j \neq i} |a_{ij}| \right\}.$$

Фигурирующие в условиях теоремы круги комплексной плоскости

$$\left\{ z \in \mathbb{C} \mid |z - a_{ii}| \leq \sum_{j \neq i} |a_{ij}| \right\}, \quad i = 1, 2, \dots, n,$$

называются *кругами Гершгорина* матрицы $A = (a_{ij})$. Можно дополнительно показать (см., к примеру, [42, 50, 95]), что если объединение кругов Гершгорина распадается на несколько связанных, но непересекающихся частей, то каждая такая часть содержит столько собственных значений матрицы, сколько кругов её составляют.

Нетрудно продемонстрировать, что теорема Гершгорина равносильна признаку Адамара неособенности матриц (Теорема 3.2.2). В самом деле, если матрица имеет диагональное преобладание, то её круги Гершгорина не захватывают начала координат комплексной плоскости, а потому в условиях теоремы Гершгорина матрица должна быть неособенной. Обратное, пусть верен признак Адамара. Если λ — собственное значение матрицы $A = (a_{ij})$, то матрица $(A - \lambda I)$ особенна и потому не может иметь диагональное преобладание. Как следствие, хотя бы для одного $i = 1, 2, \dots, n$ должно быть выполнено

$$|\lambda - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|, \quad i = 1, 2, \dots, n.$$

Этими условиями и определяются круги Гершгорина.

Пример 3.16.5 Для 2×2 -матрицы (3.10)

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix},$$

рассмотренной в Примере 3.1.3 (стр. 219), собственные значения суть $\frac{1}{2}(5 \pm \sqrt{33})$, они приблизительно равны -0.372 и 5.372 . На Рис. 3.23,

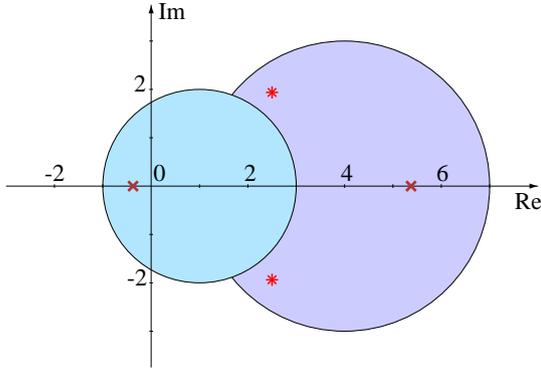


Рис. 3.23. Круги Гершгорина для матриц (3.10) и (3.11).

показывающем соответствующие матрице круги Гершгорина, эти собственные значения выделены крестиками.

Для матрицы (3.11)

$$\begin{pmatrix} 1 & 2 \\ -3 & 4 \end{pmatrix},$$

которая отличается от матрицы (3.10) лишь противоположным знаком элемента на месте (2, 1), круги Гершгорина те же. Но собственные значения у неё комплексные, равные $\frac{1}{2}(5 \pm i\sqrt{15})$, т. е. приблизительно $2.5 \pm 1.936i$. Они выделены на Рис. 3.23 звёздочками. ■

Бросается в глаза «избыточность» кругов Гершгорина, которые в качестве области локализации собственных значений очерчивают очень большую область комплексной плоскости. Это характерно для матриц с существенной внедиагональной частью. Но если недиагональные элементы матрицы малы сравнительно с диагональными, то информация, даваемая кругами Гершгорина, становится весьма точной.

3.16д Отношение Рэлея

Определение 3.16.2 Для квадратной $n \times n$ -матрицы A , вещественной или комплексной, отношением Рэлея называется функционал $\mathcal{R}(x)$, задаваемый как

$$\mathcal{R}(x) := \frac{\langle Ax, x \rangle}{\langle x, x \rangle},$$

который определён на множестве ненулевых векторов из \mathbb{R}^n или \mathbb{C}^n .

Область значений отношения Рэлея, т. е. множество

$$\{ \mathcal{R}(x) \mid x \neq 0 \},$$

называется *полем значений* матрицы A . Можно показать, что оно является выпуклым подмножеством комплексной плоскости \mathbb{C} .

Перечислим основные свойства отношения Рэлея.

Для любого скаляра α справедливо

$$\mathcal{R}(\alpha x) = \mathcal{R}(x),$$

что устанавливается непосредственной проверкой.

Если v — собственный вектор матрицы A , то $\mathcal{R}(v)$ равен собственному значению матрицы, отвечающему v . В самом деле, если обозначить это собственное значение посредством λ , то $Av = \lambda v$. По этой причине

$$\mathcal{R}(v) = \frac{\langle Av, v \rangle}{\langle v, v \rangle} = \frac{\langle \lambda v, v \rangle}{\langle v, v \rangle} = \frac{\lambda \langle v, v \rangle}{\langle v, v \rangle} = \lambda.$$

Как следствие доказанного свойства, можем заключить, что собственные числа матрицы принадлежат её полю значений.

Собственные векторы являются стационарными точками отношения Рэлея, т. е. точками зануления производной. Покажем это для вещественной симметричной матрицы, для которой отношение Рэлея рассматривается для всех ненулевых вещественных векторов:

$$\frac{\partial \mathcal{R}(x)}{\partial x_i} = \frac{\partial}{\partial x_i} \left(\frac{\langle Ax, x \rangle}{\langle x, x \rangle} \right) = \frac{2(Ax)_i \langle x, x \rangle - \langle Ax, x \rangle \cdot 2x_i}{\langle x, x \rangle^2}.$$

Если $x = v = (v_1, v_2, \dots, v_n)^\top$ — собственный вектор матрицы A , то числитель последней дроби равен $2\lambda v_i \langle v, v \rangle - \langle \lambda v, v \rangle \cdot 2v_i = 0$.

Практическое значение отношения Рэлея для вычислительных методов состоит в том, что с его помощью можно легко получить приближение к собственному значению, если известен приближённый собственный вектор матрицы.

Хотя отношение Рэлея имеет смысл и практическое значение для произвольных матриц, особую красоту и богатство содержания оно приобретает для эрмитовых (симметричных в вещественном случае) матриц.

Если A — эрмитова $n \times n$ -матрица, то, как известно,

$$A = UDU^*,$$

где $D = \text{diag} \{ \lambda_1, \lambda_2, \dots, \lambda_n \}$ — диагональная матрица с вещественными собственными значениями матрицы A по диагонали, U — некоторая унитарная $n \times n$ -матрица (ортогональная в вещественном случае). Тогда

$$\mathcal{R}(x) = \frac{\langle Ax, x \rangle}{\langle x, x \rangle} = \frac{\langle UDU^*x, x \rangle}{\langle x, x \rangle} = \frac{\langle DU^*x, U^*x \rangle}{\langle U^*x, U^*x \rangle} = \frac{\sum_{i=1}^n \lambda_i |y_i|^2}{\|y\|_2^2},$$

где $y = U^*x$. Поскольку

$$\frac{1}{\|y\|_2^2} \sum_{i=1}^n |y_i|^2 = \sum_{i=1}^n \frac{|y_i|^2}{\|y\|_2^2} = 1,$$

то для эрмитовых матриц отношение Рэля есть выпуклая комбинация, с коэффициентами $|y_i|^2 / \|y\|_2^2$, её собственных значений. В целом же из проведённых выше выкладок следует, что область значения отношения Рэля для эрмитовой матрицы — это интервал $[\lambda_{\min}, \lambda_{\max}] \subset \mathbb{R}$, коль скоро все λ_i вещественны. Кроме того, для эрмитовых матриц отношение Рэля позволяет легко находить нетривиальные границы для наименьшего собственного значения сверху и наибольшего собственного значения снизу.

В теории на основе отношения Рэля нетрудно вывести полезные оценки для собственных и сингулярных чисел матриц. В частности, из свойств отношения Рэля следует (см. подробности в [40, 50])

Теорема 3.16.3 (теорема Вейля) Пусть A и B — эрмитовы $n \times n$ -матрицы, причём $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ — собственные значения матрицы A и $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_n$ — собственные значения матрицы $\tilde{A} = A + B$. Тогда $|\lambda_i - \tilde{\lambda}_i| \leq \|B\|_2$.

Следствие. Пусть A и B — произвольные матрицы одинакового размера, причём $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ — сингулярные числа матрицы A , а $\tilde{\sigma}_1 \geq \tilde{\sigma}_2 \geq \dots \geq \tilde{\sigma}_n$ — сингулярные числа матрицы $\tilde{A} = A + B$. Тогда $|\tilde{\sigma}_i - \sigma_i| \leq \|B\|_2$.

Следствие из теоремы Вейля показывает, что сингулярные числа непрерывно зависят от элементов матрицы, и зависимость эта имеет

довольно плавный характер. В этом состоит важное отличие сингулярных чисел матрицы от её собственных чисел, которые могут изменяться в зависимости от элементов матрицы сколь угодно быстро (см. Пример 3.16.3).

3.17 Численные методы решения проблемы собственных значений

3.17а Предварительное упрощение матрицы

Естественной идеей является приведение матрицы, для которой решается проблема собственных значений, к некоторой, по-возможности, простейшей форме, для которой собственные значения и/или собственные векторы могут быть найдены проще, чем для исходной. В частности, идеальным было бы приведение матрицы к диагональной или треугольной форме, по которым собственные числа могут быть найдены непосредственно. Элементарными преобразованиями, с помощью которых может быть выполнено это приведение, в данном случае должны быть, очевидно, такие, которые сохраняют неизменным спектр матрицы, т. е. преобразования подобия матрицы $A \mapsto S^{-1}AS$. Они существенно сложнее действуют на матрицу, чем преобразования линейного комбинирования строк, которые использовались при решении систем линейных алгебраических уравнений. Невозможность полной реализации идеи упрощения матрицы следует из теоремы Абеля-Руффини, которую мы обсуждали в §3.16а: если бы это упрощение было возможным, то оно привело бы к конечному алгоритму решения алгебраических уравнений произвольной степени.

Тем не менее, в некоторых частных случаях идея предварительного упрощения матрицы для решения проблемы собственных значений, может оказаться полезной. Её наиболее популярное воплощение — это так называемая почти треугольная (хессенбергова) форма матрицы.

Определение 3.17.1 Матрица $H = (h_{ij})$ называется верхней почти треугольной или хессенберговой матрицей (в форме Хессенберга), если $h_{ij} = 0$ при $i > j + 1$.

Наглядный «портрет» хессенберговой матрицы выглядит следую-

щим образом:

$$H = \begin{pmatrix} \times & \times & \cdots & \times & \times \\ \times & \times & \cdots & \times & \times \\ & \times & \ddots & \vdots & \vdots \\ & \mathbf{0} & \ddots & \times & \times \\ & & & \times & \times \end{pmatrix}.$$

Симметричная хессенбергова матрица — это, очевидно, трёхдиагональная матрица.

Предложение 3.17.1 Любую $n \times n$ -матрицу A можно привести к ортогонально подобной хессенберговой матрице $H = QAQ^T$, где Q — произведение конечного числа отражений или вращений.

Доказательство. Рассмотрим для определённости преобразование с помощью матриц отражения.

Возьмём матрицу отражения $Q_1 = I - 2uu^T$ так, чтобы первая компонента вектора Хаусхолдера u была нулевой и при этом

$$Q_1 \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \\ \vdots \\ a_{n1} \end{pmatrix} = \begin{pmatrix} a_{11} \\ a'_{21} \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

т. е. занулялись бы элементы a_{31}, \dots, a_{n1} в первом столбце. Нетрудно видеть, что Q_1 выглядит следующим образом

$$Q_1 = \left(\begin{array}{c|cccc} 1 & 0 & \cdots & 0 & 0 \\ \hline 0 & \times & \cdots & \times & \times \\ 0 & \times & \ddots & \vdots & \vdots \\ \vdots & \vdots & \ddots & \times & \times \\ 0 & \times & \cdots & \times & \times \end{array} \right).$$

Далее, когда A умножается на такую матрицу Q_1 слева, то в ней не изменяются элементы первой строки. Когда матрица $Q_1 A$ умножается на $Q_1^T = Q_1$ справа, то в ней не изменяются элементы первого столбца.

Поэтому в матрице $Q_1 A Q_1^\top$, как и в $Q_1 A$, первый столбец имеет нули в позициях с 3-й по n -ую.

Далее выбираем матрицы отражения Q_2, Q_3, \dots, Q_{n-2} так, чтобы умножение слева на Q_i давало нули в позициях с $(i+2)$ -ой по n -ую в i -ом столбце. Эти матрицы имеют вид

$$Q_i = \left(\begin{array}{c|c} I & 0 \\ \hline 0 & \tilde{Q}_i \end{array} \right),$$

где в верхнем левом углу стоит единичная матрица размера $i \times i$, а \tilde{Q}_i — матрица отражения размера $(n-i) \times (n-i)$. При этом последующее умножения справа на $Q_i^\top = Q_i$ также не портит возникающую почти треугольную структуру результирующей матрицы. Получающаяся в итоге матрица $Q A Q^\top$ с $Q = Q_{n-2} \dots Q_1$ действительно является верхней почти треугольной. ■

3.17б Степенной метод

Если для собственных значений $\lambda_i, i = 1, 2, \dots, n$, некоторой матрицы справедливо неравенство $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots |\lambda_n|$, то λ_1 называют *доминирующим собственным значением*, а соответствующий ему собственный вектор — *доминирующим собственным вектором*. Степенной метод, описанию которого посвящён этот пункт, предназначен для решения частичной проблемы собственных значений — нахождения доминирующих собственного значения и собственного вектора матрицы.

Лежащая в его основе идея чрезвычайно проста и состоит в том, что если у матрицы A имеется собственное значение λ_1 , превосходящее по модулю все остальные собственные значения, то при действии этой матрицей на произвольный вектор $x \in \mathbb{C}^n$ направление v_1 , отвечающее этому собственному значению λ_1 будет растягиваться сильнее остальных (при $\lambda_1 > 1$) или сжиматься меньше остальных (при $\lambda_1 \leq 1$). При повторном умножении A на результат Ax предшествующего умножения эта компонента ещё более удлинится в сравнении с остальными. Повторив рассмотренную процедуру умножения достаточное количество раз, мы получим вектор, в котором полностью преобладает направление v_1 , т. е. практически будет получен приближённый собственный вектор.

В качестве приближённого собственного значения матрицы A можно при этом взять «отношение» двух последовательных векторов, порождённых нашим процессом — $x^{(k+1)} = A^{k+1}x^{(0)}$ и $x^{(k)} = A^k x^{(0)}$, $k = 0, 1, 2, \dots$. Слово «отношение» взято здесь в кавычки потому, что употреблено не вполне строго: ясно, что векторы $x^{(k+1)}$ и $x^{(k)}$ могут оказаться неколлинеарными, и тогда их «отношение» смысла иметь не будет. Возможны следующие пути решения этого вопроса:

- 1) рассматривать отношение каких-нибудь фиксированных компонент векторов $x^{(k+1)}$ и $x^{(k)}$, т. е.

$$x_i^{(k+1)} / x_i^{(k)} \quad (3.143)$$

для некоторого $i \in \{1, 2, \dots, n\}$;

- 2) рассматривать отношение проекций последовательных приближений $x^{(k+1)}$ и $x^{(k)}$ на направление, задаваемое каким-нибудь вектором $l^{(k)}$, т. е.

$$\frac{\langle x^{(k+1)}, l^{(k)} \rangle}{\langle x^{(k)}, l^{(k)} \rangle}. \quad (3.144)$$

Во втором случае мы обозначили направление проектирования через $l^{(k)}$, чтобы подчеркнуть его возможную зависимость от номера шага k . Ясно также, что это направление $l^{(k)}$ не должно быть ортогональным вектору $x^{(k)}$, чтобы не занулился знаменатель в (3.144).

Последний способ кажется более предпочтительным в вычислительном отношении, поскольку позволяет избегать капризного поведения в одной отдельно взятой компоненте вектора $x^{(k)}$, когда она может сделаться очень малой по абсолютной величине или совсем занулиться, хотя в целом вектор $x^{(k)}$ будет иметь значительную длину. Наконец, в качестве вектора, задающего направление проектирования во втором варианте, естественно взять сам $x^{(k)}$, вычисляя на каждом шаге отношение

$$\frac{\langle x^{(k+1)}, x^{(k)} \rangle}{\langle x^{(k)}, x^{(k)} \rangle}, \quad (3.145)$$

где $x^{(k)} = A^k x^{(0)}$. Нетрудно увидеть, что это выражение совпадает с отношением Рэлея для приближения $x^{(k)}$ к собственному вектору.

Для организации вычислительного алгоритма степенного метода требуется разрешить ещё два тонких момента, связанных с реализацией на ЭВМ.

Во-первых, это возможное неограниченное увеличение (при $\lambda_1 > 1$) или неограниченное уменьшение (при $\lambda_1 < 1$) норм векторов $x^{(k)}$ и $x^{(k+1)}$, участвующих в нашем процессе. Разрядная сетка современных цифровых ЭВМ, как известно, конечна и позволяет представлять числа из ограниченного диапазона. Чтобы избежать проблем, вызванных выходом за этот диапазон («переполнением» или «исчезновением порядка»), имеет смысл нормировать $x^{(k)}$. При этом наиболее удобна нормировка в евклидовой норме $\|\cdot\|_2$, так как тогда знаменатель отношения (3.145) делается равным единице.

Во-вторых, при выводе степенного метода мы неявно предполагали, что начальный вектор $x^{(0)}$ выбран так, что он имеет ненулевую проекцию на направление доминирующего собственного вектора v_1 матрицы A . В противном случае произведение любых степеней матрицы A на $x^{(0)}$ будут также иметь нулевые проекции на v_1 , и никакой дифференциации длины компонент $A^k x^{(0)}$, на которой и основывается степенной метод, не произойдёт. Это затруднение может быть преодолено с помощью какой-нибудь априорной информации о доминирующем собственном векторе матрицы. Кроме того, при практической реализации степенного метода на цифровых ЭВМ неизбежные ошибки округления, как правило, приводят к появлению ненулевых компонент в направлении v_1 , которые затем в процессе итерирования растянутся на нужную величину. Но, строго говоря, это может не происходить в некоторых исключительных случаях, и потому при ответственных вычислениях рекомендуется многократный запуск степенного метода с различными начальными векторами (так называемый мультистарт).

В псевдокоде, представленном в Табл. 3.11, $\tilde{\lambda}$ — это приближённое доминирующее собственное значение матрицы A , а $x^{(k)}$ — текущее приближение к нормированному доминирующему собственному вектору.

Теорема 3.17.1 Пусть $n \times n$ -матрица A является матрицей простой структуры (т. е. диагонализуема) и у неё имеется простое доминирующее собственное значение, которому соответствует один линейный элементарный делитель. Если начальный вектор $x^{(0)}$ не лежит в линейной оболочке $\text{lin}\{v_2, \dots, v_n\}$ собственных векторов A , которые не являются доминирующими, то степенной метод сходится.

Доказательство. При сделанных нами предположениях о матрице A она может быть представлена в виде

$$A = VDV^{-1},$$

Таблица 3.11. Степенной метод для нахождения доминирующего собственного значения матрицы

```

k ← 0;
выбираем вектор x(0) ≠ 0;
нормируем x(0) ← x(0) / ||x(0)||2;
DO WHILE ( метод не сошёлся )
    y(k+1) ← Ax(k);
    λ̃ ← ⟨y(k+1), x(k)⟩;
    x(k+1) ← y(k+1) / ||y(k+1)||2;
    k ← k + 1;
END DO

```

где $D = \text{diag} \{ \lambda_1, \lambda_2, \dots, \lambda_n \}$ — диагональная матрица с собственными значениями $\lambda_1, \lambda_2, \dots, \lambda_n$ по диагонали, а V — матрица, осуществляющая преобразование подобия, причём без ограничения общности можно считать, что λ_1 — доминирующее собственное значение A . Матрица V составлена из собственных векторов v_i матрицы A как из столбцов:

$$V = (v_1 \ v_2 \ \dots \ v_n) = \left(\begin{array}{c|c|c|c} (v_1)_1 & (v_2)_1 & \dots & (v_n)_1 \\ (v_1)_2 & (v_2)_2 & \dots & (v_n)_2 \\ \vdots & \vdots & \ddots & \vdots \\ (v_1)_n & (v_2)_n & \dots & (v_n)_n \end{array} \right),$$

где через $(v_i)_j$ обозначена j -ая компонента i -го собственного вектора

матрицы A . При этом можно считать, что $\|v_i\|_2 = 1$. Следовательно,

$$\begin{aligned} A^k x^{(0)} &= (VDV^{-1})^k x^{(0)} = \underbrace{(VDV^{-1})(VDV^{-1}) \dots (VDV^{-1})}_{k \text{ раз}} x^{(0)} \\ &= VD(V^{-1}V)D(V^{-1}V) \dots (V^{-1}V)DV^{-1}x^{(0)} \\ &= VD^k V^{-1}x^{(0)} = VD^k z \\ &= V \begin{pmatrix} \lambda_1^k z_1 \\ \lambda_2^k z_2 \\ \vdots \\ \lambda_n^k z_n \end{pmatrix} = (\lambda_1^k z_1) V \begin{pmatrix} 1 \\ (\lambda_2/\lambda_1)^k (z_2/z_1) \\ \vdots \\ (\lambda_n/\lambda_1)^k (z_n/z_1) \end{pmatrix}, \end{aligned}$$

где обозначено $z = V^{-1}x^{(0)}$. Необходимое условие последнего преобразования этой цепочки — $z_1 \neq 0$ — выполнено потому, что в условиях теоремы вектор $x^{(0)} = Vz$ должен иметь ненулевую первую компоненту при разложении по базису из собственных векторов A , т. е. столбцов матрицы V .

Коль скоро λ_1 — доминирующее собственное значение матрицы A , т. е.

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|,$$

то все частные $\lambda_2/\lambda_1, \lambda_3/\lambda_1, \dots, \lambda_n/\lambda_1$ по модулю меньше единицы, и потому при $k \rightarrow \infty$ вектор

$$\begin{pmatrix} 1 \\ (\lambda_2/\lambda_1)^k (z_2/z_1) \\ \vdots \\ (\lambda_n/\lambda_1)^k (z_n/z_1) \end{pmatrix} \tag{3.146}$$

сходится к вектору $(1, 0, 0, \dots, 0)^T$. Соответственно, произведение

$$V \begin{pmatrix} 1 \\ (\lambda_2/\lambda_1)^k (z_2/z_1) \\ \vdots \\ (\lambda_n/\lambda_1)^k (z_n/z_1) \end{pmatrix}$$

сходится к первому столбцу матрицы V , т.е. к собственному вектору, отвечающему λ_1 . Вектор $x^{(k)}$, который отличается от $A^k x^{(0)}$ лишь нормировкой, сходится к собственному вектору v_1 , а величина $\tilde{\lambda} = \langle y^{(k+1)}, x^{(k)} \rangle$ сходится к $\langle Av_1, v_1 \rangle = \langle \lambda_1 v_1, v_1 \rangle = \lambda_1$. ■

Из проведённых выше выкладок следует, что быстрота сходимости степенного метода определяется отношениями $|\lambda_i/\lambda_1|$, $i = 2, 3, \dots, n$, — знаменателями геометрических прогрессий, стоящих в качестве элементов вектора (3.146). Фактически, решающее значение имеет наибольшее из этих отношений, т.е. $|\lambda_2/\lambda_1|$, зависящее от того, насколько модуль доминирующего собственного значения отделён от модуля остальной части спектра. Чем больше эта отделённость, тем быстрее сходимость степенного метода.

Пример 3.17.1 Для матрицы (3.10)

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix},$$

при вычислениях с двойной точностью степенной метод с начальным вектором $x^{(0)} = (1, 1)^T$ за 7 итераций даёт семь верных знаков доминирующего собственного значения $\frac{1}{2}(5 + \sqrt{33}) \approx 5.3722813$. Детальная картина сходимости показана в следующей табличке:

Номер итерации	Приближение к собственному значению
1	5.0
2	5.3448276
3	5.3739445
4	5.3721649
5	5.3722894
6	5.3722808
7	5.3722814

Быстрая сходимость объясняется малостью величины $|\lambda_2/\lambda_1|$, которая, как мы могли видеть в Примере 3.1.3, для рассматриваемой матрицы равна всего лишь 0.069.

Для матрицы (3.11)

$$\begin{pmatrix} 1 & 2 \\ -3 & 4 \end{pmatrix},$$

при тех же исходных условиях степенной метод порождает последовательность значений $\tilde{\lambda}$, которая случайно колеблется от примерно 0.9 до 4 и очевидным образом не имеет предела. Причина — наличие у матрицы двух одинаковых по абсолютной величине комплексно-сопряжённых собственных значений $2.5 \pm 1.936i$ (см. Пример 3.1.3). ■

Отметим, что для симметричных (эрмитовых) положительно определённых матриц в степенном методе в качестве приближения к доминирующему собственному значению можно брать отношение

$$\frac{\|x^{(k+1)}\|_2}{\|x^{(k)}\|_2}, \quad x^{(k+1)} = Ax^{(k)}$$

(см. [77]).

Наконец, необходимое замечание о сходимости степенного метода в комплексном случае. Так как комплексные числа описываются парами вещественных чисел, то комплексные одномерные инвариантные пространства матрицы имеют вещественную размерность 2. Даже будучи нормированными, векторы из такого подпространства могут отличаться на скалярный множитель $e^{i\varphi}$ для какого-то аргумента φ , так что если не принять специальных мер, то в степенном методе видимой стабилизации координатных представлений комплексных собственных векторов может не наблюдаться. Тем не менее, о факте сходимости или расходимости можно при этом судить по стабилизации приближения к собственному значению. Либо кроме нормировки собственных векторов следует предусмотреть ещё приведение их к такой форме, в которой координатные представления будут определяться более «жёстко», например, требованием, чтобы первая компонента вектора была бы чисто вещественной.

Пример 3.17.2 Рассмотрим работу степенного метода в применении к матрице

$$\begin{pmatrix} 1 & 2i \\ 3 & 4i \end{pmatrix},$$

имеющей собственные значения

$$\lambda_1 = -0.4308405 - 0.1485958i,$$

$$\lambda_2 = 1.4308405 + 4.1485958i.$$

Доминирующим собственным значением здесь является λ_2 .

Начав итерирование с вектора $x^{(0)} = (1, 1)^\top$, уже через 7 итераций мы получим 6 правильных десятичных знаков в вещественной и мнимой частях собственного значения. Но вот в порождаемых алгоритмом нормированных векторах $x^{(k)}$ —

$$x^{(9)} = \begin{pmatrix} -0.01132 - 0.43223 i \\ -0.11659 - 0.89413 i \end{pmatrix},$$

$$x^{(10)} = \begin{pmatrix} 0.40491 - 0.15163 i \\ 0.80725 - 0.40175 i \end{pmatrix},$$

$$x^{(11)} = \begin{pmatrix} 0.27536 + 0.33335 i \\ 0.64300 + 0.63215 i \end{pmatrix},$$

$$x^{(12)} = \begin{pmatrix} 0.22535 + 0.36900 i \\ -0.38795 + 0.81397 i \end{pmatrix},$$

и так далее — нелегко «невооружённым глазом» узнать один и тот же собственный вектор, который «крутится» в одномерном комплексном инвариантном подпространстве. Но если поделить все получающиеся векторы на их первую компоненту, то получим один и тот же результат

$$\begin{pmatrix} 1. \\ 2.07430 - 0.21542 i \end{pmatrix},$$

и теперь уже налицо факт сходимости собственных векторов. ■

Как ведёт себя степенной метод в случае, когда матрица A не является диагонализуемой? Полный анализ ситуации можно найти, например, в книгах [42, 44]. Наиболее неблагоприятен при этом случай, когда доминирующее собственное значение находится в жордановой клетке размера два и более. Теоретически степенной метод всё таки сходится к этому собственному значению, но уже медленнее любой геометрической прогрессии.

Пример 3.17.3 Рассмотрим работу степенного метода в применении к матрице

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix},$$

т. е. к жордановой 2×2 -клетке с собственным значением 1.

Запустив степенной метод из начального вектора $x^{(0)} = (1, 1)^T$, будем иметь следующее

Номер итерации	Приближение к собственному значению
1	1.5
3	1.3
10	1.0990099
30	1.0332963
100	1.009999
300	1.0033333
1000	1.001

То есть, для получения n верных десятичных знаков собственного значения приходится делать примерно 10^{n-1} итераций, что, конечно же, непомерно много. При увеличении размера жордановой клетки сходимость степенного метода делается ещё более медленной. ■

3.17в Обратные степенные итерации

Обратными степенными итерациями для матрицы A называют описанный в прошлом параграфе степенной метод, применённый к обратной матрице A^{-1} , в котором вычисляется отношение результатов предыдущей итерации к последующей, т. е. обратная к (3.143) или (3.144) величина. Явное нахождение обратной матрицы A^{-1} при этом не требуется, так как в степенном методе используется лишь результат $x^{(k+1)}$ её умножения на вектор $x^{(k)}$ очередного приближения, а это, как известно (см., в частности, §3.13), эквивалентно решению системы линейных уравнений $Ax^{(k+1)} = x^{(k)}$.

Так как собственные значения матриц A и A^{-1} взаимно обратны, то обратные степенные итерации будут сходиться к наименьшему по абсолютной величине собственному значению A и соответствующему собственному вектору.

Чтобы в отношении

$$\frac{\langle x^{(k)}, l^{(k)} \rangle}{\langle x^{(k+1)}, l^{(k)} \rangle},$$

которое необходимо вычислять в обратных степенных итерациях, знаменатель не занулялся, удобно брать $l^{(k)} = x^{(k+1)}$. Тогда очередным

Таблица 3.12. Обратные степенные итерации для нахождения наименьшего по модулю собственного значения матрицы A

```

 $k \leftarrow 0$ ;
выбираем вектор  $x^{(0)} \neq 0$ ;
нормируем  $x^{(0)} \leftarrow x^{(0)} / \|x^{(0)}\|_2$ ;
DO WHILE ( метод не сошёлся )
    найти  $y^{(k+1)}$  из системы  $Ay^{(k+1)} = x^{(k)}$ ;
 $\tilde{\lambda} \leftarrow \langle x^{(k)}, y^{(k+1)} \rangle / \langle y^{(k+1)}, y^{(k+1)} \rangle$ ;
 $x^{(k+1)} \leftarrow y^{(k+1)} / \|y^{(k+1)}\|_2$ ;
 $k \leftarrow k + 1$ ;
END DO

```

приближением к наименьшему по модулю собственному значению матрицы A является

$$\frac{\langle x^{(k)}, x^{(k+1)} \rangle}{\langle x^{(k+1)}, x^{(k+1)} \rangle},$$

где $Ax^{(k+1)} = x^{(k)}$. Псевдокод получающегося метода представлен в Табл. 3.12.

Практическая реализация решения системы линейных уравнений (5-я строка псевдокода) может быть сделана достаточно эффективной, если предварительно выполнить LU- или QR-разложение матрицы A , а затем на каждом шаге метода использовать формулы (3.59) или (3.77).

Пример 3.17.4 Рассмотрим работу обратных степенных итераций для знакомой нам матрицы

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix},$$

собственные значения которой суть $\frac{1}{2}(5 \pm \sqrt{33})$, приблизительно равные -0.372 и 5.372 .

Запустив обратные степенные итерации из начального вектора $x^{(0)} = (1, 1)^T$, за 7 итераций получим 7 верных значащих цифр наименьшего по модулю собственного числа 0.3722813. Скорость сходимости здесь получается такой же, как в Примере 3.14.1 для доминирующего собственного значения этой матрицы, что неудивительно ввиду одинакового значения знаменателя геометрической прогрессии λ_2/λ_1 . ■

Обратные степенные итерации особенно эффективны в случае, когда имеется хорошее приближение к собственному значению и требуется найти соответствующий собственный вектор.

3.17г Сдвиги спектра

Сдвигом матрицы называют прибавление к ней скалярной матрицы, т. е. матрицы, пропорциональной единичной матрице, так что вместо матрицы A мы получаем матрицу $A + \vartheta I$. Если $\lambda_i(A)$ — собственные значения матрицы A , то для любого комплексного числа ϑ собственными значениями матрицы $A + \vartheta I$ являются числа $\lambda_i(A) + \vartheta$, тогда как собственные векторы остаются неизменными. Цель сдвига — преобразование спектра матрицы для того, чтобы улучшить работу тех или иных алгоритмов решения проблемы собственных значений.

Если, к примеру, у матрицы A наибольшими по абсолютной величине были два собственных значения -2 и 2 , то прямое применение к ней степенного метода не приведёт к успеху. Но у матрицы $A + I$ эти собственные значения перейдут в -1 и 3 , второе собственное число станет наибольшим по модулю, и теперь уже единственным. Соответственно, степенной метод делается применимым к новой матрице.

Пример 3.17.5 Для матрицы (3.11)

$$\begin{pmatrix} 1 & 2 \\ -3 & 4 \end{pmatrix},$$

как было отмечено в Примере 3.17б, простейший степенной метод расходится из-за существования двух наибольших по абсолютной величине собственных значений.

Но если сдвинуть эту матрицу на $2i$, то её спектр (см. Рис. 3.23) поднимется «вверх», абсолютные величины собственных значений перестанут совпадать, и степенной метод окажется применимым.

Степенные итерации для «сдвинутой» матрицы

$$\begin{pmatrix} 1 + 2i & 2 \\ -3 & 4 + 2i \end{pmatrix} \tag{3.147}$$

довольно быстро сходятся к наибольшему по модулю собственному значению $\frac{5}{2} + (2 + \frac{1}{2}\sqrt{15})i \approx 2.5 + 3.9364917i$. Детальная картина сходимости при вычислениях с двойной точностью и начальным вектором $x^{(0)} = (1, 1)^T$ показана в следующей табличке:

Номер итерации	Приближение к собственному значению
1	$2.0 + 2.0i$
3	$2.0413223 + 4.3140496i$
5	$2.7022202 + 3.9372711i$
10	$2.5004558 + 3.945456i$
20	$2.4999928 + 3.9364755i$

В данном случае для матрицы (3.147) имеем $|\lambda_2/\lambda_1| \approx 0.536$. ■

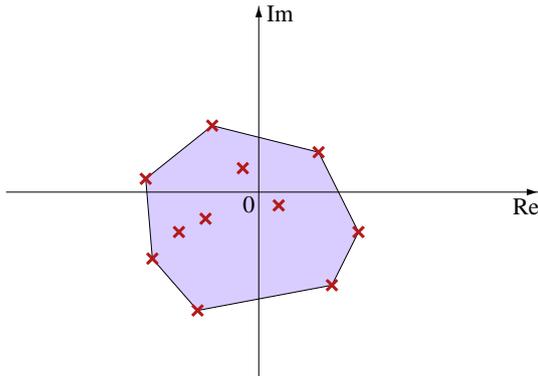


Рис. 3.24. С помощью подходящих сдвигов любую крайнюю точку спектра можно сделать наибольшей по модулю.

Поскольку спектр симметричной (эрмитовой) матрицы лежит на вещественной оси, то к таким матрицам имеет смысл применять вещественные сдвиги. В частности, при этом для симметричных веще-

ственных матриц алгоритмы будут реализовываться в более простой вещественной арифметике.

С помощью сдвигов матрицы можно любое её собственное значение, которое является крайней точкой выпуклой оболочки спектра, сделать наибольшим по модулю, обеспечив, таким образом, сходимость к нему итераций степенного метода. Но как добиться сходимости к другим собственным значениям, которые лежат «внутри» спектра, а не «с краю»? Здесь на помощь приходят обратные степенные итерации.

Обратные степенные итерации сходятся к ближайшей к нулю точке спектра матрицы, и такой точкой с помощью подходящего сдвига может быть сделано любое собственное число. В этом — важное преимущество сдвигов для обратных степенных итераций.

Другое важное следствие сдвигов — изменение отношения $|\lambda_2/\lambda_1|$, величина которого влияет на скорость сходимости степенного метода. Обычно с помощью подходящего выбора величины сдвига ϑ можно добиться того, чтобы

$$\left| \frac{\lambda_2 + \vartheta}{\lambda_1 + \vartheta} \right|$$

было меньше, чем $|\lambda_2/\lambda_1|$, ускорив тем самым степенные итерации. Совершенно аналогичный эффект оказывает удачный выбор сдвига на отношение $|\lambda_n/\lambda_{n-1}|$, которое определяет скорость сходимости обратных степенных итераций.

3.17д Метод Якоби для решения симметричной проблемы собственных значений

В этом параграфе мы рассмотрим численный метод для решения симметричной проблемы собственных значений, т. е. для вычисления собственных чисел и собственных векторов симметричных матриц. Он был впервые применён К.Г. Якоби в 1846 году к конкретной 7×7 -матрице, а затем был забыт на целое столетие и вновь переоткрыт лишь после Второй мировой войны, когда началось бурное развитие вычислительной математики.

Идея метода Якоби состоит в том, чтобы подходящими преобразованиями подобия от шага к шагу уменьшать норму внедиагональной части матрицы. Получающиеся при этом матрицы имеют тот же спектр, что и исходная матрица, но будут стремиться к диагональной матрице с собственными значениями на главной диагонали. Инструментом реализации этого плана выступают элементарные ортогональные матрицы

вращений, рассмотренные в §3.7д. Почему именно ортогональные матрицы и почему вращений? Ответ на эти вопросы станет ясен позднее при анализе работы алгоритма.

Итак, положим $A^{(0)} := A$. Если матрица $A^{(k)}$, $k = 0, 1, 2, \dots$, уже вычислена, то подберём матрицу вращений $G(p, q, \theta)$ вида (3.81) таким образом, чтобы сделать нулями пару внедиагональных элементов в позициях (p, q) и (q, p) в матрице $A^{(k+1)} := G(p, q, \theta)^\top A^{(k)} G(p, q, \theta)$. Желая достичь этой цели, мы должны добиться выполнения равенства

$$\begin{pmatrix} a_{pp}^{(k+1)} & a_{pq}^{(k+1)} \\ a_{qp}^{(k+1)} & a_{qq}^{(k+1)} \end{pmatrix} = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}^\top \begin{pmatrix} a_{pp}^{(k)} & a_{pq}^{(k)} \\ a_{qp}^{(k)} & a_{qq}^{(k)} \end{pmatrix} \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \\ = \begin{pmatrix} \times & 0 \\ 0 & \times \end{pmatrix},$$

где, как обычно, посредством « \times » обозначены какие-то элементы, конкретное значение которых несущественно. Строго говоря, в результате рассматриваемого преобразования подобия в матрице $A^{(k)}$ изменятся и другие элементы, находящиеся в строках и столбцах с номерами p и q . Этот эффект будет проанализирован ниже в Предложении 3.17.3.

Опуская индексы, обозначающие номер итерации и приняв сокращённые обозначения $c = \cos \theta$, $s = \sin \theta$, получим

$$\begin{pmatrix} \times & 0 \\ 0 & \times \end{pmatrix} = \begin{pmatrix} a_{pp}c^2 + a_{qq}s^2 + 2sca_{pq} & sc(a_{qq} - a_{pp}) + a_{pq}(c^2 - s^2) \\ sc(a_{qq} - a_{pp}) + a_{pq}(c^2 - s^2) & a_{pp}s^2 + a_{qq}c^2 - 2sca_{pq} \end{pmatrix}$$

Приравнивание внедиагональных элементов нулю даёт

$$\frac{a_{pp} - a_{qq}}{a_{pq}} = \frac{c^2 - s^2}{sc}.$$

Поделив обе части этой пропорции пополам, воспользуемся тригонометрическими формулами двойных углов

$$\frac{a_{pp} - a_{qq}}{2a_{pq}} = \frac{c^2 - s^2}{2sc} = \frac{\cos(2\theta)}{\sin(2\theta)} = \frac{1}{\operatorname{tg}(2\theta)} =: \tau.$$

Положим $t := \sin \theta / \cos \theta = \operatorname{tg} \theta$. Вспоминая далее тригонометрическую формулу для тангенса двойного угла

$$\operatorname{tg}(2\theta) = \frac{2 \operatorname{tg} \theta}{1 - \operatorname{tg}^2 \theta},$$

мы можем прийти к выводу, что t является корнем квадратного уравнения

$$t^2 + 2\tau t - 1 = 0$$

с положительным дискриминантом ($4\tau^2 + 4$), которое, следовательно, всегда имеет вещественные корни. Отсюда находится сначала t :

$$t = -\tau \pm \sqrt{\tau^2 + 1},$$

причём из двух корней мы берём наименьший по абсолютной величине. Он равен

$$\begin{aligned} t &= -\tau + \operatorname{sgn} \tau \cdot \sqrt{\tau^2 + 1}, & \text{если } \tau \neq 0, \\ t &= \pm 1, & \text{если } \tau = 0, \end{aligned}$$

и первую формулу для улучшения численной устойчивости лучше записать в виде, освобождённом от вычитания близких чисел:

$$\begin{aligned} t &= \frac{(-\tau + \operatorname{sgn} \tau \cdot \sqrt{\tau^2 + 1})(\tau + \operatorname{sgn} \tau \cdot \sqrt{\tau^2 + 1})}{(\tau + \operatorname{sgn} \tau \cdot \sqrt{\tau^2 + 1})} \\ &= \frac{1}{(\tau + \operatorname{sgn} \tau \cdot \sqrt{\tau^2 + 1})} \quad \text{при } \tau \neq 0. \end{aligned}$$

Затем на основе известных тригонометрических формул, выражающих косинус и синус через тангенс, находим c и s :

$$c = \frac{1}{\sqrt{t^2 + 1}}, \quad s = t \cdot c.$$

Займёмся теперь обоснованием сходимости метода Якоби для решения симметричной проблемы собственных значений.

Предложение 3.17.2 Фробениусова норма матрицы A , т. е.

$$\|A\|_F = \left(\sum_{i,j=1}^n a_{ij}^2 \right)^{1/2},$$

не изменяется при умножениях на ортогональные матрицы слева или справа.

Доказательство. Напомним, что *следом матрицы* $A = (a_{ij})$, обозначаемым $\text{tr } A$, называется сумма всех её диагональных элементов:

$$\text{tr } A = \sum_{i=1}^n a_{ii}.$$

Нетрудно проверить, что привлечение понятия следа позволяет переписать определение фробениусовой нормы матрицы таким образом

$$\|A\|_F = \left(\sum_{j=1}^n \left(\sum_{i=1}^n a_{ij} a_{ij} \right) \right)^{1/2} = (\text{tr}(A^T A))^{1/2}.$$

Следовательно, для любой ортогональной матрицы Q справедливо

$$\begin{aligned} \|QA\|_F &= \left(\text{tr}((QA)^T(QA)) \right)^{1/2} \\ &= \left(\text{tr}(A^T Q^T Q A) \right)^{1/2} = (\text{tr}(A^T A))^{1/2} = \|A\|_F. \end{aligned}$$

Для доказательства аналогичного соотношения с умножением на ортогональную матрицу справа заметим, что фробениусова норма не меняется при транспонировании матрицы. Следовательно,

$$\|AQ\|_F = \|(Q^T A^T)^T\|_F = \|Q^T A^T\|_F = \|A^T\|_F = \|A\|_F,$$

что завершает доказательство Предложения. ■

Следствие. Фробениусова норма матрицы не меняется при ортогональных преобразованиях подобия.

Для более точного описания меры близости матриц $A^{(k)}$, которые порождаются конструируемым нами методом, к диагональной матрице введём величину

$$ND(A) = \left(\sum_{j \neq i} a_{ij}^2 \right)^{1/2}$$

— фробениусову норму внедиагональной части матрицы. Ясно, что матрица A диагональна тогда и только тогда, когда $ND(A) = 0$.

Предложение 3.17.3 Пусть преобразование подобия матрицы A с помощью матрицы вращений G таково, что в матрице $B = G^T A G$ зануляются элементы в позициях (p, q) и (q, p) . Тогда

$$ND^2(B) = ND^2(A) - 2a_{pq}^2. \tag{3.148}$$

Итак, в сравнении с матрицей A в матрице B изменились элементы строк и столбцов с номерами p и q , но фробениусова норма недиагональной части изменилась при этом так, как будто кроме зануления элементов a_{pq} и a_{qp} ничего не произошло.

Доказательство. Для 2×2 -подматрицы

$$\begin{pmatrix} a_{pp} & a_{pq} \\ a_{qp} & a_{qq} \end{pmatrix}$$

из матрицы A и соответствующей ей 2×2 -подматрицы

$$\begin{pmatrix} b_{pp} & 0 \\ 0 & b_{qq} \end{pmatrix}$$

в матрице B справедливо соотношение

$$a_{pp}^2 + a_{qq}^2 + 2a_{pq}^2 = b_{pp}^2 + b_{qq}^2,$$

так как ортогональным преобразованием подобия фробениусова норма матрицы не изменяется. Но, кроме того, $\|A\|_F^2 = \|B\|_F^2$, и потому

$$\begin{aligned} ND^2(B) &= \|B\|_F^2 - \sum_{i=1}^n b_{ii}^2 \\ &= \|A\|_F^2 - \left(\sum_{i=1}^n a_{ii}^2 - (a_{pp}^2 + a_{qq}^2) + (b_{pp}^2 + b_{qq}^2) \right) \\ &= ND^2(A) - 2a_{pq}^2, \end{aligned}$$

поскольку на диагонали у матрицы A изменились только два элемента — a_{pp} и a_{qq} . ■

Таблица 3.13. Метод Якоби для вычисления собственных значений симметричной матрицы

<p>Вход</p> <p>Симметричная матрица A.</p> <p>Допуск ϵ на норму внедиагональных элементов.</p>
<p>Выход</p> <p>Матрица, на диагонали которой стоят приближения к собственным значениям A.</p>
<p>Алгоритм</p> <p>DO WHILE ($ND(A) > \epsilon$)</p> <p style="padding-left: 40px;">выбрать ненулевой внедиагональный элемент a_{pq} в A;</p> <p style="padding-left: 40px;">обнулить a_{pq} и a_{qp} преобразованием подобия с матрицей вращения $G(p, q, \theta)$;</p> <p>END DO</p>

Теперь можно ответить на вопрос о том, почему в методе Якоби для преобразований подобия применяются именно ортогональные матрицы. Как следует из результатов Предложений 3.17.2 и 3.17.3, умножение на ортогональные матрицы обладает замечательным свойством сохранения фробениусовой нормы матрицы и, как следствие, «перекачивания» её величины с внедиагональных элементов на диагональ в результате специально подобранных цепочек таких умножений. При других преобразованиях подобия добиться этого было бы едва ли возможно.

Итак, всё готово для организации итерационного процесса приведения симметричной матрицы к диагональному виду, при котором внедиагональные элементы последовательно подавляются. Как уже отмечалось, занулённые на каком-то шаге алгоритма элементы могут впоследствии вновь сделаться ненулевыми. Но результат Предложения 3.17.3 показывает, что норма внедиагональной части матрицы при этом всё равно монотонно уменьшается.

Различные способы выбора ненулевых внедиагональных элементов, подлежащих обнулению, приводят к различным практическим версиям метода Якоби.

Выбор наибольшего по модулю внедиагонального элемента — наилучшее для отдельно взятого шага алгоритма решение. Но поиск этого элемента имеет трудоёмкость $n(n - 1)/2$, что может оказаться весьма дорогостоящим, особенно для матриц больших размеров. Преобразование подобия с матрицей вращений обходится всего в $O(n)$ операций!

Чаще применяют циклический обход столбцов (или строк) матрицы, и наибольший по модулю элемент берут в пределах рассматриваемого столбца (строки).

Наконец, ещё одна популярная версия — это так называемый «барьерный метод Якоби», в котором назначают величину «барьера» на значение модуля внедиагональных элементов матрицы, и алгоритм обнуляет все элементы, модуль которых превосходит этот барьер. Затем барьер понижается, процесс обнуления повторяется заново, и так до тех пор, пока не будет достигнута требуемая точность.

К 70-м годам прошлого века, когда было разработано немало эффективных численных методов для решения симметричной проблемы собственных значений, стало казаться, что метод Якоби устарел и будет вытеснен из широкой вычислительной практики (см., к примеру, рассуждения в [78]). Дальнейшее развитие не подтвердило эти пессимистичные прогнозы. Выяснилось, что метод Якоби почти не имеет конкурентов по точности нахождения малых собственных значений, тогда как методы, основанные на трёхдиагонализации исходной матрицы, могут терять точность (соответствующие примеры приведены в [13]). Кроме того, метод Якоби оказался хорошо распараллеливаемым, т. е. подходящим для расчётов на современных многопроцессорных ЭВМ.

3.17e Базовый QR-алгоритм

QR-алгоритм, изложению которого посвящён этот параграф, является одним из наиболее эффективных численных методов для решения полной проблемы собственных значений. Он был изобретён независимо В.Н. Кублановской (1960 год) и Дж. Фрэнсисом (1961 год). Публикация В.Н. Кублановской появилась раньше²⁹, а Дж. Фрэнсис более пол-

²⁹Упомянув о вкладе В.Н. Кублановской в изобретение QR-алгоритма, обычно ссылаются на её статью 1961 года в «Журнале вычислительной математики и математической физики» [75]. Но первое сообщение о QR-алгоритме было опубликовано

но развил практичную версию QR-алгоритма.

QR-алгоритм является наиболее успешным представителем большого семейства родственных методов решения полной проблемы собственных значений, основанных на разложении исходной матрицы на простейшие. QR-алгоритму предшествовал LR-алгоритм Рутисхаузера. На практике применяются также ортогональный степенной метод, предложенный В.В. Воеводиным, и различные другие близкие вычислительные процессы.

Вспомним теорему о QR-разложении (Теорема 3.7.1, стр. 302): всякая квадратная матрица представима в виде произведения ортогональной и правой (верхней) треугольной матриц. Ранее в нашем курсе мы уже обсуждали конструктивные способы выполнения этого разложения — с помощью матриц отражения Хаусхолдера, а также с помощью матриц вращений. Следовательно, далее можно считать, что QR-разложение выполнимо и основывать на этом факте свои построения.

Вычислительная схема базового QR-алгоритма для решения проблемы собственных значений представлена в Табл. 3.14: мы разлагаем матрицу $A^{(k)}$, полученную на k -м шаге алгоритма, $k = 0, 1, 2, \dots$, на ортогональный $Q^{(k)}$ и правый треугольный $R^{(k)}$ сомножители и далее, поменяв их местами, умножаем друг на друга, образуя следующее приближение $A^{(k+1)}$.

Таблица 3.14. QR-алгоритм для нахождения собственных значений матрицы A

```

 $k \leftarrow 0;$ 
 $A^{(0)} \leftarrow A;$ 
DO WHILE ( метод не сошёлся )
    вычислить QR-разложение  $A^{(k)} = Q^{(k)} R^{(k)}$ ;
     $A^{(k+1)} \leftarrow R^{(k)} Q^{(k)}$ ;
     $k \leftarrow k + 1;$ 
END DO

```

ею раньше — в Дополнении к изданию 1960 года книги [44].

Прежде всего отметим, что поскольку

$$A^{(k+1)} = R^{(k)}Q^{(k)} = (Q^{(k)})^\top (Q^{(k)}R^{(k)})Q^{(k)} = (Q^{(k)})^\top A^{(k)}Q^{(k)},$$

то все матрицы $A^{(k)}$, $k = 0, 1, 2, \dots$, ортогонально подобны друг другу и исходной матрице A . Как следствие, собственные значения всех матриц $A^{(k)}$ совпадают с собственными значениями A . Результат о сходимости QR-алгоритма неформальным образом может быть резюмирован в следующем виде: если A — неособенная вещественная матрица, то последовательность порождаемых QR-алгоритмом матриц $A^{(k)}$ сходится «по форме» к верхней блочно-треугольной матрице.

Это означает, что предельная матрица, к которой сходится QR-алгоритм, является верхней треугольной либо верхней блочно-треугольной, причём размеры диагональных блоков зависят, во-первых, от типа собственных значений матрицы (кратности и принадлежности вещественной оси \mathbb{R}), и, во-вторых, от того, в вещественной или комплексной арифметике выполняется QR-алгоритм.

Если алгоритм выполняется в вещественной (комплексной) арифметике и все собственные значения матрицы вещественны (комплексны) и различны по модулю, то предельная матрица — верхняя треугольная. Если алгоритм выполняется в вещественной (комплексной) арифметике и некоторое собственное значение матрицы вещественно (комплексно) и имеет кратность p , то в предельной матрице ему соответствует диагональный блок размера $p \times p$. Если алгоритм выполняется для вещественной матрицы в вещественной арифметике, то простым комплексно-сопряжённым собственным значениям (они имеют равные модули) отвечают диагональные 2×2 -блоки в предельной матрице. Наконец, если некоторое комплексное собственное значение вещественной матрицы имеет кратность p , так что ему соответствует ещё такое же комплексно-сопряжённое собственное значение кратности p , то при выполнении QR-алгоритма в вещественной арифметике предельная матрица получит диагональный блок размера $2p \times 2p$.

Пример 3.17.6 Проиллюстрируем работу QR-алгоритма на примере матрицы

$$\begin{pmatrix} 1 & -2 & 3 \\ 4 & 5 & -6 \\ -7 & 8 & 9 \end{pmatrix}, \quad (3.149)$$

имеющей собственные значения

$$\begin{aligned} &2.7584148 \\ &6.1207926 \pm 8.0478897i \end{aligned}$$

Читатель может провести на компьютере этот увлекательный эксперимент самостоятельно, воспользовавшись системами Scilab, MATLAB или им подобными: все они имеют встроенную процедуру для QR-разложения матриц.³⁰ ■

Пример 3.17.7 Для ортогональной матрицы

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad (3.150)$$

QR-разложением является произведение её самой на единичную матрицу. Поэтому в результате одного шага QR-алгоритма мы снова получим исходную матрицу, которая, следовательно, и будет пределом итераций. В то же время, матрица (3.150) имеет собственные значения, равные ± 1 , так что в данном случае QR-алгоритм не работает. ■

3.17ж Модификации QR-алгоритма

Представленная в Табл. 3.14 версия QR-алгоритма на практике обычно снабжается рядом модификаций, которые существенно повышают её эффективность. Главными из этих модификаций являются

- 1) сдвиги матрицы, рассмотренные нами в §3.17г, и
- 2) предварительное приведение матрицы к специальной верхней почти треугольной форме.

Можно показать (см. теорию в книгах [13, 41]), что, аналогично степенному методу, сдвиги также помогают ускорению QR-алгоритма. Но в QR-алгоритме их традиционно организуют способом, представленным в Табл. 3.15.

Особенность организации сдвигов в этом псевдокоде — присутствие обратных сдвигов (в строке 6 алгоритма) сразу же вслед за прямыми (в 5-й строке). Из-за этого в получающемся алгоритме последовательно

³⁰В Scilab'е и MATLAB'е она так и называется — `qr`.

Таблица 3.15. QR-алгоритм со сдвигами для нахождения собственных значений матрицы A

```

 $k \leftarrow 0;$ 
 $A^{(0)} \leftarrow A;$ 
DO WHILE ( метод не сошёлся )
    выбрать сдвиг  $\vartheta_k$  вблизи собственного значения  $A$ ;
    вычислить QR-разложение  $A^{(k)} - \vartheta_k I = Q^{(k)} R^{(k)}$ ;
     $A^{(k+1)} \leftarrow R^{(k)} Q^{(k)} + \vartheta_k I$ ;
     $k \leftarrow k + 1$ ;
END DO

```

вычисляемые матрицы $A^{(k)}$ и $A^{(k+1)}$ ортогонально подобны, совершенно так же, как и в исходной версии QR-алгоритма:

$$\begin{aligned}
 A^{(k+1)} &= R^{(k)} Q^{(k)} + \vartheta_k I = (Q^{(k)})^\top Q^{(k)} R^{(k)} Q^{(k)} + \vartheta_k (Q^{(k)})^\top Q^{(k)} \\
 &= (Q^{(k)})^\top (Q^{(k)} R^{(k)} + \vartheta_k I) Q^{(k)} = (Q^{(k)})^\top A^{(k)} Q^{(k)}.
 \end{aligned}$$

То есть, представленная организация сдвигов позволила сделать их в одно и то же время локальными и динамическими по характеру.

Пример 3.17.8 Проиллюстрируем работу QR-алгоритма со сдвигами на знакомой нам матрице (3.149)

$$\begin{pmatrix} 1 & -2 & 3 \\ 4 & 5 & -6 \\ -7 & 8 & 9 \end{pmatrix}$$

из предыдущего примера

■

Предложение 3.17.4 Матрица, имеющая хессенбергову форму, сохраняет эту форму при выполнении с ней QR-алгоритма.

Доказательство. При QR-разложении хессенберговой матрицы в качестве ортогонального сомножителя Q для матрицы $(A^{(k)} - \vartheta I)$ получается также хессенбергова матрица, так как j -ый столбец в Q есть линейная комбинация первых j столбцов матрицы $(A^{(k)} - \vartheta I)$. В свою очередь, матрица RQ — произведение после перестановки сомножителей — также получается хессенберговой. Добавление диагонального слагаемого ϑI не изменяет верхней почти треугольной формы матрицы. ■

Смысл предварительного приведения к хессенберговой форме заключается в следующем. Хотя это приведение матрицы требует $O(n^3)$ операций, дальнейшее выполнение одной итерации QR-алгоритма с хессенберговой формой будет теперь стоить всего $O(n^2)$ операций, так что общая трудоёмкость QR-алгоритма составит $O(n^3)$. Для исходной версии QR-алгоритма, которая оперирует с плотно заполненной матрицей, трудоёмкость равна $O(n^4)$, поскольку на каждой итерации алгоритма выполнение QR-разложения требует $O(n^3)$ операций.

3.18 Численные методы нахождения сингулярных чисел и векторов

Сингулярные числа зависят от элементов матрицы существенно более плавным образом, нежели собственные числа. Мы могли видеть это из следствия из теоремы Вейля (теорема 3.16.3). На эту тему существует ещё один известный результат

Теорема 3.18.1 (теорема Виландта-Хофмана) Пусть A и B — эрмитовы $n \times n$ -матрицы, причём $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ — собственные значения матрицы A и $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_n$ — собственные значения матрицы $\tilde{A} = A + B$. Тогда

$$\left(\sum_{i=1}^n (\tilde{\lambda}_i - \lambda_i)^2 \right)^{1/2} \leq \|B\|_F,$$

где $\|\cdot\|_F$ — фробениусова норма матрицы.

Доказательство можно найти в [41, 42]

Простейшие методы нахождения сингулярных чисел матриц основаны на том, что они являются собственными числами матриц $A^T A$ и AA^T ($A^* A$ и AA^* в комплексном случае).

Литература к главе 3

Основная

- [1] Бахвалов Н.С., Жидков Н.П., Кобельков Г.М. *Численные методы*. – Москва: Бином, 2003, а также другие издания этой книги.
- [2] Бахвалов Н.С., Корнев А.А., Чижонков Е.В. *Численные методы. Решения задач и упражнения*. – Москва: Дрофа, 2008.
- [3] Березин И.С., Жидков Н.П. *Методы вычислений. Т. 1–2*. – Москва: Наука, 1966.
- [4] Вержбицкий В.М. *Численные методы. Части 1–2*. – Москва: «Оникс 21 век», 2005.
- [5] Воеводин В.В. *Вычислительные основы линейной алгебры*. – Москва: Наука, 1977.
- [6] Воеводин В.В. *Линейная алгебра*. – Москва: Наука, 1980.
- [7] Воеводин В.В., Воеводин Вл.В. *Энциклопедия линейной алгебры. Электронная система ЛИНЕАЛ*. – Санкт-Петербург: БХВ-Петербург, 2006.
- [8] Волков Е.А. *Численные методы*. – Москва: Наука, 1987.
- [9] Гантмахер Ф.Р. *Теория матриц*. – Москва: Наука, 1988.
- [10] Глазман И.М., Лювич Ю.И. *Конечномерный линейный анализ*. – Москва: Наука, 1969.
- [11] Голуб Дж., ван Лоун Ч. *Матричные вычисления*. – Москва: Мир, 1999.
- [12] Демидович Б.П., Марон А.А. *Основы вычислительной математики*. – Москва: Наука, 1970.
- [13] Деммель Дж. *Вычислительная линейная алгебра*. – Москва: Мир, 2001.
- [14] Зорич В.А. *Математический анализ. Т. 1*. – Москва: Наука, 1981. *Т. 2*. – Москва: Наука, 1984, а также более поздние издания.
- [15] Икрамов Х.Д. *Численные методы для симметричных линейных систем*. – Москва: Наука, 1988.
- [16] Ильин В.П. *Методы и технологии конечных элементов*. – Новосибирск: Издательство ИВМиМГ СО РАН, 2007.
- [17] Ильин В.П., Кузнецов Ю.И. *Трёхдиагональные матрицы и их приложения*. – Москва: Наука, 1985.
- [18] Канторович Л.В., Акилов Г.П. *Функциональный анализ*. – Москва: Наука, 1984.
- [19] Като Т. *Теория возмущений линейных операторов*. – Москва: Мир, 1972.

- [20] Коллатц Л. *Функциональный анализ и вычислительная математика*. – Москва: Мир, 1969.
- [21] Коновалов А.Н. *Введение в вычислительные методы линейной алгебры*. – Новосибирск: Наука, 1993.
- [22] Кострикин А.Н. *Введение в алгебру. Часть 1. Основы алгебры*. – Москва: Физматлит, 2004.
- [23] Кострикин А.Н. *Введение в алгебру. Часть 2. Линейная алгебра*. – Москва: Физматлит, 2000.
- [24] Красносельский М.А., Крейн С.Г. Итеративный процесс с минимальными невязками // *Математический Сборник*. – 1952. – Т. 31 (73), №2. – С. 315–334.
- [25] Крылов В.И., Бобков В.В., Монастырный П.И. *Вычислительные методы. Т. 1–2*. – Москва: Наука, 1976.
- [26] Ланкастер П. *Теория матриц*. – Москва: Наука, 1978.
- [27] Лоусон Ч., Хенсон Р. *Численное решение задач методом наименьших квадратов*. – Москва: Наука, 1986.
- [28] Марчук Г.И., Кузнецов Ю.А. *Итерационные методы и квадратичные функционалы*. – Новосибирск: Наука, 1972.
- [29] *Матрицы и квадратичные формы. Основные понятия. Терминология* / Академия Наук СССР. Комитет научно-технической терминологии. – Москва: Наука, 1990. – (Сборники научно-нормативной терминологии; Вып. 112).
- [30] Мацокин А.М. *Численный анализ. Вычислительные методы линейной алгебры. Конспекты лекций для преподавания в III семестре ММФ НГУ*. – Новосибирск: НГУ, 2009–2010.
- [31] Миньков С.Л., Миньков Л.Л. *Основы численных методов*. – Томск: Издательство научно-технической литературы, 2005.
- [32] Мысовских И.П. *Лекции по методам вычислений*. – Санкт-Петербург: Издательство Санкт-Петербургского университета, 1998.
- [33] Ортега Дж. *Введение в параллельные и векторные методы решения линейных систем*. – Москва: Мир, 1991.
- [34] Островский А.М. *Решение уравнений и систем уравнений*. – Москва: Издательство иностранной литературы, 1963.
- [35] Прасолов В.В. *Задачи и теоремы линейной алгебры*. – Москва: Наука-Физматлит, 1996.
- [36] Райс Дж. *Матричные вычисления и математическое обеспечение*. – Москва: Мир, 1984.
- [37] Самарский А.А., Гулин А.В. *Численные методы*. – Москва: Наука, 1989.
- [38] Стренг Г. *Линейная алгебра и её применения*. – Москва: Мир, 1980.
- [39] Тихонов А.Н., Арсенин В.Я. *Методы решения некорректных задач*. – Москва: Наука, 1974.
- [40] Тыртышников Е.Е. *Матричный анализ и линейная алгебра*. – Москва: Физматлит, 2007.

- [41] Тыртышников Е.Е. *Методы численного анализа*. – Москва: Академия, 2007.
- [42] Уилкинсон Дж. *Алгебраическая проблема собственных значений*. – Москва: Наука, 1970.
- [43] Уоткинс Д. *Основы матричных вычислений*. – Москва: «Бином. Лаборатория знаний», 2009.
- [44] Фаддеев Д.К., Фаддеева В.Н. *Вычислительные методы линейной алгебры*. – Москва–Ленинград: Физматлит, 1960 (первое издание) и 1963 (второе издание).
- [45] Федоренко Р.П. Итерационные методы решения разностных эллиптических уравнений // *Успехи Математических Наук*. – 1973. – Т. 28, вып. 2 (170). – С. 121–182.
- [46] Форсайт Дж.Э. Что представляют собой релаксационные методы? // *Современная математика для инженеров под ред. Э.Ф.Беккенбаха*. – Москва: Издательство иностранной литературы, 1958. – С. 418–440.
- [47] Форсайт Дж., Молер К. *Численное решение систем линейных алгебраических уравнений*. – Москва: Мир, 1969.
- [48] Хаусдорф Ф. *Теория множеств*. – Москва: УРСС Эдиториал, 2007.
- [49] Хейгеман Л., Янг Д. *Прикладные итерационные методы*. – Москва: Мир, 1986.
- [50] Хорн Р., Джонсон Ч. *Матричный анализ*. – Москва: Мир, 1989.
- [51] Шилов Г.Е. *Математический анализ. Конечномерные линейные пространства*. – Москва: Наука, 1969.
- [52] Шилов Г.Е. *Математический анализ. Функции одного переменного. Часть 3*. – Москва: Наука, 1970.
- [53] Авертн О. *Precise numerical methods using C++*. – San Diego: Academic Press, 1998.
- [54] BECKERMANN B. The condition number of real Vandermonde, Krylov and positive definite Hankel matrices // *Numerische Mathematik*. – 2000. – Vol. 85, No. 4. – P. 553–577.
- [55] KELLEY C.T. *Iterative methods for linear and nonlinear equations*. – Philadelphia: SIAM, 1995.
- [56] SAAD Y. *Iterative methods for sparse linear systems*. – Philadelphia: SIAM, 2003.
- [57] Scilab — The Free Platform for Numerical Computation. <http://www.scilab.org>
- [58] TEMPLE G. The general theory of relaxation methods applied to linear systems // *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences*. – 1939. – Vol. 169, No. 939. – P. 476–500.
- [59] TREFETHEN L.N., BAU D. III *Numerical linear algebra*. – Philadelphia: SIAM, 1997.

Дополнительная

- [60] АЛЕКСАНДРОВ П.С. *Введение в теорию множеств и общую топологию.* – Санкт-Петербург: Лань, 2010.
- [61] АЛЕКСЕЕВ В.Б. *Теорема Абеля в задачах и решениях.* – Москва: Московский Центр непрерывного математического образования, 2001.
- [62] АЛЕКСЕЕВ Е.Р., ЧЕСНОКОВА О.В., РУДЧЕНКО Е.А. *Scilab. Решение инженерных и математических задач.* – Москва: Alt Linux – Бином, 2008.
- [63] АЛЕФЕЛЬД Г., ХЕРЦБЕРГЕР Ю. *Введение в интервальные вычисления.* – Москва: Мир, 1987.
- [64] ВОЕВОДИН В.В., КУЗНЕЦОВ Ю.А. *Матрицы и вычисления.* – Москва: Наука, 1984.
- [65] ГОДУНОВ С.К., АНТОНОВ А.Г., КИРИЛЮК О.Г., КОСТИН В.И. *Гарантированная точность решения систем линейных уравнений в евклидовых пространствах.* – Новосибирск: Наука, 1988 и 1992.
- [66] ГОДУНОВ С.К. *Современные аспекты линейной алгебры.* – Новосибирск: Научная книга, 1997.
- [67] ГОРБАЧЕНКО В.И. *Вычислительная линейная алгебра с примерами на MATLAB.* – Санкт-Петербург: «БХВ-Петербург», 2011.
- [68] ДЖОРДЖ А., ЛЮ ДЖ. *Численное решение больших разреженных систем уравнений.* – Москва: Мир, 1984.
- [69] ДРОБЫШЕВИЧ В.И., ДЫМНИКОВ В.П., РИВИН Г.С. *Задачи по вычислительной математике.* – Москва: Наука, 1980.
- [70] ЗЕЛЬДОВИЧ Я.Б., МЫШКИС А.Д. *Элементы прикладной математики.* – Москва: Наука, 1967.
- [71] ИКРАМОВ Х.Д. *Численные методы линейной алгебры.* – Москва: Знание, 1987.
- [72] ИКРАМОВ Х.Д. *Численное решение матричных уравнений.* – Москва: Наука, 1984.
- [73] КАЛИТКИН Н.Н. *Численные методы.* – Москва: Наука, 1978.
- [74] КРЫЛОВ А.Н. *Лекции о приближённых вычислениях.* – Москва: ГИТТЛ, 1954, а также более ранние издания.
- [75] КУБЛАНОВСКАЯ В.Н. О некоторых алгоритмах для решения полной проблемы собственных значений // *Журнал вычисл. матем. и мат. физики.* – 1961. – Т. 1, № 4. – С. 555–570.
- [76] КУЗНЕЦОВ Ю.А. Метод сопряжённых градиентов, его обобщения и применения // *Вычислительные процессы и системы.* – Москва: Наука, 1983 – Вып. 1. – С. 267–301.
- [77] МАРЧУК Г.И. *Методы вычислительной математики.* – Москва: Наука, 1989.
- [78] ПАРЛЕТТ Б. *Симметричная проблема собственных значений. Численные методы.* – Москва: Мир, 1983.
- [79] ПАРОДИ М. *Локализация характеристических чисел матриц и её применения.* – Москва: Издательство иностранной литературы, 1960.

- [80] САМАРСКИЙ А.А., НИКОЛАЕВ Е.С. *Методы решения сеточных уравнений.* – Москва: Наука, 1978.
- [81] ФАДДЕЕВА В.Н. *Вычислительные методы линейной алгебры.* – Москва–Ленинград: Гостехиздат, 1950.
- [82] ФЛЭНАГАН Д., МАЦУМОТО Ю. *Язык программирования Ruby.* – Санкт-Петербург: Питер, 2011.
- [83] ХАЛМОШ П. *Конечномерные векторные пространства.* – Москва: ГИФМЛ, 1963.
- [84] ШАРЫЙ С.П. *Конечномерный интервальный анализ.* – Электронная книга, 2012 (см. <http://www.nsc.ru/interval/Library/InteBooks>)
- [85] ЯНЕНКО Н.Н. *Метод дробных шагов решения многомерных задач математической физики.* – Новосибирск: Наука, 1967.
- [86] BAUER F.L., FIKE C.T. Norms and exclusion theorems // *Numerische Mathematik.* – 1960. – Vol. 2. – P. 137–141.
- [87] ЕСКАРТ С., YOUNG G. The approximation of one matrix by another of lower rank // *Psychometrika.* – 1936. – Vol. 1. – P. 211–218.
- [88] GREGORY R.T., KARNEY D.L. *A collection of matrices for testing computational algorithms.* – Hantington, New York: Robert E. Krieger Publishing Company, 1978.
- [89] HOTELLING H. Analysis of a complex of statistical variables into principal components // *J. Educ. Psych.* – 1933 – Vol. 24. – Part I: pp. 417-441, Part II: pp. 498-520.
- [90] KREINOVICH V., LAKEYEV A.V., NOSKOV S.I. Approximate linear algebra is intractable // *Linear Algebra and its Applications.* – 1996. – Vol. 232. – P. 45–54.
- [91] KREINOVICH V., LAKEYEV A.V., ROHN J., KAHL P. *Computational complexity and feasibility of data processing and interval computations.* – Dordrecht: Kluwer, 1997.
- [92] MOLER C. Professor SVD // *The MathWorks News & Notes.* – October 2006. – P. 26–29.
- [93] MOORE R.E., KEARFOTT R.B., CLOUD M. *Introduction to interval analysis.* – Philadelphia: SIAM, 2009.
- [94] SCHULZ G. Iterative Berechnung der reziproken Matrix // *Z. Angew. Math. Mech.* – 1933. – Bd. 13 (1). – S. 57–59.
- [95] STOER J., BULIRSCH R. *Introduction to numerical analysis.* – Berlin-Heidelberg-New York: Springer-Verlag, 1993.
- [96] VARGA R.S. *Matrix iterative analysis.* – Berlin, Heidelberg, New York: Springer Verlag, 2000, 2010.
- [97] WILF H.S. *Finite sections of some classical inequalities.* – Heidelberg: Springer, 1970.
- [98] TODD J. The condition number of the finite segment of the Hilbert matrix // *National Bureau of Standards, Applied Mathematics Series.* – 1954. – Vol. 39. – P. 109–116.

Глава 4

Решение нелинейных уравнений и их систем

4.1 Введение

В этой главе рассматривается задача решения системы уравнений

$$\begin{cases} F_1(x_1, x_2, \dots, x_n) = 0, \\ F_2(x_1, x_2, \dots, x_n) = 0, \\ \vdots \quad \ddots \quad \vdots \\ F_n(x_1, x_2, \dots, x_n) = 0, \end{cases} \quad (4.1)$$

над полем вещественных чисел \mathbb{R} , или, кратко,

$$F(x) = 0, \quad (4.2)$$

где $x = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ — вектор неизвестных переменных,

$F_i(x)$, $i = 1, 2, \dots, n$, — вещественнозначные функции,

$F(x) = (F_1(x), F_2(x), \dots, F_n(x))^T$ — вектор-столбец функций F_i .

Для переменных x_1, x_2, \dots, x_n нужно найти набор значений $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$, называемый *решением системы*, который обращает в равенства все уравнения системы (4.1)–(4.2). В некоторых случаях желательно найти все такие возможные наборы, т. е. все решения системы, иногда достаточно какого-то одного. В случае, когда система уравнений (4.1)–(4.2)

не имеет решений, нередко требуется предоставить обоснование этого заключения или даже его вывод, и им может быть протокол работы программы для ЭВМ и т. п.

Наряду с задачами, рассмотренными в Главе 2, то есть интерполяции и приближениями функций, вычислением интегралов, задача решения уравнений и систем уравнений является одной из классических задач вычислительной математики.

Всюду далее мы предполагаем, что функции $F_i(x)$ по меньшей мере непрерывны, а количество уравнений в системе (4.1)–(4.2) совпадает с количеством неизвестных переменных. Помимо записи систем уравнений в каноническом виде (4.1)–(4.2) часто встречаются и другие формы их представления, например,

$$G(x) = H(x) \quad (4.3)$$

с какими-то функциями G, H . Чрезвычайно важным частным случаем этой формы является *рекуррентный вид* системы уравнений (уравнения),

$$x = G(x). \quad (4.4)$$

в котором неизвестная переменная выражена через саму себя. В этом случае решение системы уравнений (или уравнения) есть *неподвижная точка* отображения G , т. е. такой элемент области определения G , который переводится этим отображением сам в себя. Кроме того, рекуррентный вид уравнения или системы хорош тем, что позволяет довольно просто организовать итерационный процесс для нахождения решения, что мы могли видеть в Главе 3.

Как правило, системы уравнений различного вида могут быть приведены друг к другу равносильными преобразованиями. В частности, несложно установить связь решений уравнений и систем уравнений вида (4.1)–(4.2) с неподвижными точками отображений, т. е. с решениями уравнений в рекуррентном виде (4.4). Ясно, что

$$F(x) = 0 \quad \Longleftrightarrow \quad x = x - \Lambda F(x),$$

где Λ — ненулевой скаляр в одномерном случае или же неособенная $n \times n$ -матрица в случае вектор-функции F . Поэтому решение уравнения

$$F(x) = 0$$

является неподвижной точкой отображения

$$G(x) = x - \Lambda F(x).$$

Если неизвестная x не является конечномерным вектором, а отображения F и G имеют весьма общую природу, то математические свойства уравнений (4.2) и (4.4) могут существенно различаться, так что при этом формы записи (4.2) и (4.4), строго говоря, не вполне равносильны друг другу. По этой причине для их обозначения часто употребляют отдельные термины — *уравнение первого рода* и *уравнение второго рода* соответственно.

Обращаясь к решению нелинейных уравнений и их систем, мы обнаруживаем себя в гораздо более сложных условиях, нежели при решении систем линейных алгебраических уравнений (3.43)–(3.44). Стройная и весьма полная теория разрешимости систем линейных уравнений, базирующаяся на классических результатах линейной алгебры, обеспечивала в необходимых нам случаях уверенность в существовании решения систем линейных уравнений и его единственности. Для нелинейных уравнений столь общей и простой теории не существует. Напротив, нелинейные уравнения и их системы имеют в качестве общего признака лишь отрицание линейности, т.е. то, что все они «не линейны», и потому отличаются огромным разнообразием. Из общих нелинейных уравнений и систем уравнений принято выделять *алгебраические* уравнения и системы уравнений, в которых функции $F_i(x)$ являются алгебраическим полиномами относительно неизвестных переменных x_1, x_2, \dots, x_n .

4.2 Вычислительно-корректные задачи

4.2а Предварительные сведения и определения

Напомним общеизвестный факт: на вычислительных машинах (как электронных, так и механических, как цифровых, так и аналоговых) мы можем выполнять, как правило, лишь приближённые вычисления над полем вещественных чисел \mathbb{R} . Для цифровых вычислительных машин этот вывод следует из того, что они являются дискретными и конечными устройствами, так что и ввод вещественных чисел в такую вычислительную машину и выполнение с ними различных арифметических операций сопровождаются неизбежными ошибками, вызванными конечным характером представления чисел, конечностью исполнительных устройств и т.п. Для аналоговых вычислительных машин данные также не могут быть введены абсолютно точно, и процесс вычислений также не абсолютно точен. Потенциально все отмеченные погрешности

могут быть сделаны сколь угодно малыми, но в принципе избавиться от них не представляется возможным. Получается, что реально

- мы решаем на вычислительной машине не исходную математическую задачу, а более или менее близкую к ней,
- сам процесс решения на ЭВМ отличается от своего идеального математического прообраза, т. е. от результатов вычислений в \mathbb{R} или \mathbb{C} по тем формулам, которые его задают.

Возникновение и бурное развитие компьютерной алгебры с её «безошибочными» вычислениями едва ли опровергает высказанный выше тезис, так как исходные постановки задач для систем символьных преобразований требуют *точную* представимость входных данных, которые поэтому подразумеваются целыми или, на худой конец, рациональными с произвольной длиной числителя и знаменателя (см. [2]), а все преобразования над ними не выводят за пределы поля рациональных чисел.

Как следствие, в условиях приближённого представления входных числовых данных и приближённого характера вычислений над полем вещественных чисел \mathbb{R} мы в принципе можем решать лишь те постановки задач, ответы которых «не слишком резко» меняются при изменении входных данных, т. е. устойчивы по отношению к возмущениям в этих начальных данных. Для этого, по крайней мере, должна иметь место непрерывная зависимость решения от входных данных.

Для формализации высказанных выше соображений нам необходимо точнее определить ряд понятий.

Под *массовой задачей* [13] будем понимать некоторый общий вопрос, формулировка которого содержит несколько свободных переменных — *параметров* — могущих принимать значения в пределах предписанных им множеств. В целом массовая задача Π определяется

- 1) указанием её входных данных, т. е. общим списком всех *параметров* с областями их определения,
- 2) формулировкой тех свойств, которым должен удовлетворять *ответ*, т. е. решение задачи.

Индивидуальная задача I получается из массовой задачи Π путём присваивания всем параметрам задачи Π каких-то конкретных значений. Наконец, *разрешающим отображением* задачи Π мы называем отображение, сопоставляющее каждому набору входных данных-параметров

ответ соответствующей индивидуальной задачи (см. §1.3). Станем говорить, что массовая математическая задача является *вычислительно корректной*, если её разрешающее отображение $\mathcal{P} \rightarrow \mathcal{A}$ из множества входных данных \mathcal{P} во множество \mathcal{A} ответов задачи непрерывно относительно некоторых топологий на \mathcal{P} и \mathcal{A} , определяемых содержательным смыслом задачи.

Те задачи, ответы на которые неустойчивы по отношению к возмущениям входных данных, могут решаться на ЭВМ с конечной разрядной сеткой лишь опосредованно, после проведения мероприятий, необходимых для защиты от этой неустойчивости или её нейтрализации.

Конечно, скорость изменения решения в зависимости от изменений входных данных может быть столь большой, что эта зависимость, даже будучи непрерывной и сколь угодно гладкой, становится похожей на разрывную. Это мы могли видеть в §3.16в для собственных значений некоторых матриц, которые являются «практически разрывными» функциями элементов матрицы. Но определением вычислительно корректной задачи выделяются те задачи, для которых хотя бы в принципе возможно добиться сколь угодно точного приближения к идеальному математическому ответу, например, увеличением количества значащих цифр при вычислениях и т. п.

Пример 4.2.1 Задача решения систем линейных уравнений $Ax = b$ с неособенной квадратной матрицей A является вычислительно-корректной. Если топология на пространстве \mathbb{R}^n её решений задается обычным евклидовым расстоянием и подобным же традиционным образом задётся расстояние между векторами правой части и матрицами, то существуют хорошо известные неравенства (см. §3.5а), оценивающие сверху границы изменения решений x через изменения элементов матрицы A , правой части b и число обусловленности матрицы A . ■

Пример 4.2.2 Вычисление ранга матрицы — вычислительно некорректная задача. Дело в том, что в основе понятия ранга лежит линейная зависимость строк или столбцов матрицы, т. е. свойство, которое нарушается при сколь угодно малых возмущениях матрицы. ■

Разрывная зависимость решения от входных данных задачи может возникать вследствие присутствия в алгоритме вычисления функции условных операторов вида IF ... THEN ... ELSE, приводящих к ветвле-

нию. Такова хорошо известная функция знака числа

$$\operatorname{sgn} x = \begin{cases} -1, & \text{если } x < 0, \\ 0, & \text{если } x = 0, \\ 1, & \text{если } x > 0. \end{cases}$$

Аналогична функция модуля числа $|x|$, с которой в обычных и внешне простых выражениях могут быть замаскированы разрывы и ветвления. Например, таково частное $\sin x / |x|$, которое ведёт себя в окрестности нуля примерно как $\operatorname{sgn} x$.

Для систем нелинейных уравнений, могущих иметь неединственное решение, топологию на множестве ответов \mathcal{A} нужно задавать уже каким-либо расстоянием между множествами, например, с помощью так называемой *хаусдорфовой метрики* [8]. Напомним её определение.

Если задано метрическое пространство с метрикой ρ , то *расстоянием* точки a до множества X называется величина $\rho(a, X)$, определяемая как $\inf_{x \in X} \rho(a, x)$. *Хаусдорфовым расстоянием* между компактными множествами X и Y называют величину

$$\rho(X, Y) = \max \left\{ \max_{x \in X} \rho(x, Y), \max_{y \in Y} \rho(y, X) \right\}.$$

При этом $\rho(X, Y) = +\infty$, если $X = \emptyset$ или $Y = \emptyset$. Введённая таким образом величина действительно обладает всеми свойствами расстояния и может быть использована для задания топологии на пространствах решений тех задач, ответы к которым неединственны, т. е. являются целыми множествами.

4.26 Задача решения уравнений не является вычислительно-корректной

Уже простейшие примеры показывают, что задача решения уравнений и систем уравнений не является вычислительно-корректной. Например, квадратное уравнение

$$x^2 + px + q = 0 \tag{4.5}$$

для

$$p^2 = 4q \tag{4.6}$$

имеет лишь одно решение $x = -p/2$. Но при любых сколь угодно малых возмущениях коэффициента p и свободного члена q , нарушающих

равенство (4.6), уравнение (4.5) теряет это единственное решение или же приобретает ещё одно (см. Рис. 4.1).

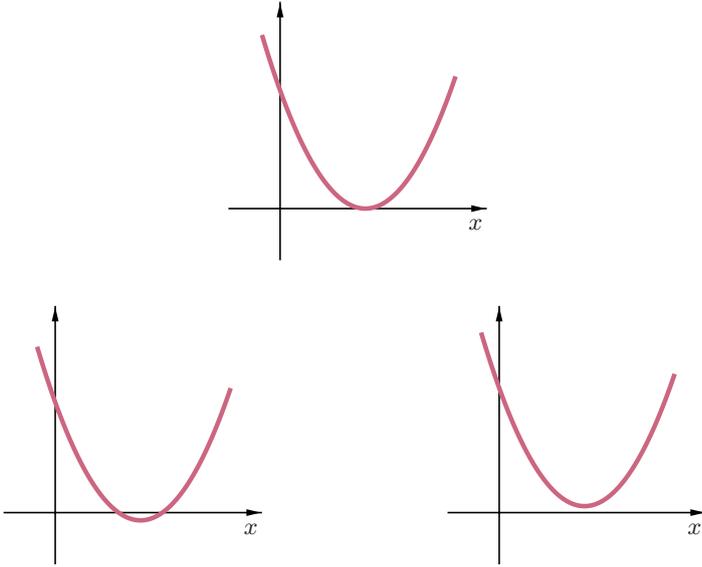


Рис. 4.1. Неустойчивая зависимость решений уравнения (4.5)–(4.6) от сколь угодно малых шевелений его коэффициентов.

Аналогичным образом ведёт себя решение двумерной системы уравнений, эквивалентной (4.5),

$$\begin{cases} x + y = r, \\ xy = s \end{cases}$$

при $s = r^2/4$. При этом раздвоение решения не является большим грехом, коль скоро мы можем рассматривать хаусдорфово расстояние между целостными множествами решений. Но вот исчезновение единственного решения, при котором расстояние между множествами решений скачком меняется до $+\infty$ — это чрезвычайное событие, однозначно указывающее на разрывность разрешающего отображения.

Как видим, математическую постановку задачи нахождения решений уравнений нужно «исправить», заменив какой-нибудь вычислительно-корректной постановкой задачи. Приступая к поиску ответа на

этот математический вопрос, отметим, прежде всего, что с точки зрения практических приложений задачи, которые мы обычно формулируем в виде решения уравнений или систем уравнений, традиционно выписывая соотношение

$$F(x) = 0 \quad (4.2)$$

и ему подобные, имеют весьма различную природу. Это и будет отправной точкой нашей ревизии постановки задачи решения уравнений и систем уравнений.

4.2в ε -решения уравнений

В ряде практических задач пользователю требуется не точное равенство некоторого выражения нулю, а лишь его «исчезающая малость» в сравнении с каким-то а priori установленным порогом. С аналогичной точки зрения часто имеет смысл рассматривать соотношения вида (4.3) или (4.4), выражающие равенство двух каких-то выражений.

Таковы, например, в большинстве физических, химических и других естественнонаучных расчётов уравнения материального баланса, вытекающие из закона сохранения массы и закона сохранения заряда. Точное равенство левой и правой частей уравнения здесь неявным образом и не требуется, так как погрешность этого равенства всегда ограничена снизу естественными пределами делимости материи. В самом деле, масса молекулы, масса и размеры атома, заряд элементарной частицы и т. п. величины, с точностью до которых имеет смысл рассматривать конкретные уравнения баланса — все они имеют вполне конечные (хотя и весьма малые) значения.

Например, не имеет смысла требовать, чтобы закон сохранения заряда выполнялся с погрешностью, меньшей чем величина элементарного электрического заряда (т. е. заряда электрона, равного $1.6 \cdot 10^{-19}$ Кл). Также бессмысленно требовать, чтобы погрешность изготовления или подгонки деталей оптических систем была существенно меньшей длины световой волны (от $4 \cdot 10^{-7}$ м до $7.6 \cdot 10^{-7}$ м в зависимости от цвета). А что касается температуры, то при обычных земных условиях определение её с точностью, превосходящей 0.001 градуса, вообще проблематично в силу принципиальных соображений. Наконец, ограниченная точность, с которой известны абсолютно все физические константы¹,

¹В лучшем случае относительная погрешность известных на сегодняшний день значений физических констант равна 10^{-10} , см. [58].

также воздвигает границы для требований равенства в физических соотношениях.

Совершенно аналогична ситуация с экономическими балансами, как в стоимостном выражении, так и в натуральном: требовать, чтобы они выполнялись с погрешностью, меньшей, чем одна копейка (наименьшая денежная величина) или чем единица неделимого товара (телевизор, автомобиль и т. п.) просто бессмысленно.

Во всех вышеприведённых примерах под решением уравнения понимается значение переменной, которое доставляет левой и правой частям уравнения пренебрежимо отличающиеся значения. В применении к уравнениям вида (4.2) соответствующая формулировка выглядит следующим образом:

Для заданных отображения $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ и $\varepsilon > 0$ найти значения неизвестной переменной x , такие что $F(x) \approx 0$ с абсолютной погрешностью ε , т. е. $\|F(x)\| < \varepsilon$.

Решением этой задачи является целое множество точек, которые мы будем называть ε -решениями или почти решениями, если порог этой пренебрежимой малости не оговорён явно или несуществен.

Нетрудно понять, что условием $\|F(x)\| < \varepsilon$ задаётся открытое множество, если отображение F непрерывно. Любая точка из этого множества устойчива к малым возмущениям исходных данных, а задача «о нахождении почти решений» является вычислительно-корректной.

Как уже отмечалось выше, в некоторых задачах система уравнений более естественно записывается не как (4.2), а в виде (4.3)

$$G(x) = H(x),$$

и требуется обеспечить с относительной погрешностью ε равенство её левой и правой частей:

Для заданных отображений $G, H : \mathbb{R}^n \rightarrow \mathbb{R}^n$ и $\varepsilon > 0$ найти значения неизвестной переменной x , такие что $G(x) \approx H(x)$ с относительной погрешностью ε , т. е.

$$\frac{\|G(x) - H(x)\|}{\max\{\|G(x)\|, \|H(x)\|\}} < \varepsilon .$$

Решения этой задачи мы также будем называть ε -решениями системы уравнений вида (4.3).

Математические понятия, определения которых привлекают малый допуск ε , не являются чем-то экзотическим. Таковы, к примеру, ε -энтропия множеств в метрических пространствах, ε -субдифференциал функции, ε -оптимальные решения задач оптимизации и т. п. Одним из частных случаев ε -решений являются точки ε -спектра матрицы, предложенные для обобщения традиционного понятия собственного значения матрицы [11, 47, 59]. Говорят, что точка z на комплексной плоскости принадлежит ε -спектру матрицы A , если существует комплексный вектор v единичной длины, такой что $\|(A - zI)v\| \leq \varepsilon$, где $\|\cdot\|$ — какая-то векторная норма. Иными словами, при условии $\|v\| = 1$ здесь рассматривается приближённое «с точностью до ε » равенство $Av = zv$.

4.2г Недостаточность ε -решений

Но есть и принципиально другой тип задач, которые образно могут быть названы задачами «об определении перехода через нуль» и не сводятся к нахождению ε -решений. Таковы задачи, в которых требуется гарантированно отследить переход функции к значениям противоположного знака (или, более общо, переход через некоторое критическое значение). При этом, в частности, в любой окрестности решения должны присутствовать как положительные значения функции, так и её отрицательные значения, тогда как в задачах нахождения «почти решений» это условие может и не выполняться.

Рассмотрим следующую ситуацию, для анализа которой достаточно знание элементарной физики. Пусть кирпич лежит на опоре (см. Рис. 4.2), и мы потихоньку сдвигаем его к краю. Когда он упадёт? Для

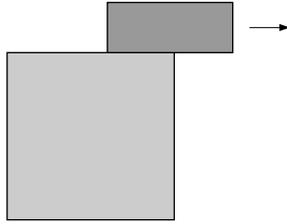


Рис. 4.2. Когда кирпич упадёт с подставки?

ответа на этот вопрос приравнивают момент силы тяжести, действующей на свисающую часть кирпича, и момент силы тяжести, действующей на ту часть, которая лежит на опоре.

Но в случае точного их равенства кирпич ещё не упадёт! Эта ситуация называется неустойчивым равновесием, но в отсутствие каких-либо воздействий на кирпич он не будет падать, а зависнет на грани опоры. Для падения кирпича именно нужен его переход чуть дальше этого положения неустойчивого равновесия. В частности, ε -решения здесь не годятся по существу дела.

Другой пример. Фазовый переход в физической системе (плавление, кристаллизация и т. п.) — типичная задача такого сорта, так как в процессе фазового перехода температура системы не меняется. Если мы хотим узнать, прошёл ли фазовый переход полностью, то нужно зафиксировать момент достижения множества состояний, лежащего по другую сторону от границы раздела различных состояний!

Ещё один пример. Рассмотрим систему линейных дифференциальных уравнений с постоянными коэффициентами

$$\frac{dx}{dt} = Ax, \quad (4.7)$$

матрица которой $A = A(\theta)$ зависит от параметра θ (возможно, векторного). Пусть при некотором начальном значении $\theta = \theta_0$ собственные значения $\lambda(A)$ матрицы A имеют отрицательные вещественные части, так что все решения системы (4.7) устойчивы по Ляпунову (и даже асимптотически устойчивы). При каких значениях параметра θ рассматриваемая система делается неустойчивой?

Традиционно отвечают на этот вопрос следующим образом. Срыв устойчивости в системе (4.7) произойдет при $\operatorname{Re} \lambda(A(\theta)) = 0$ для какого-

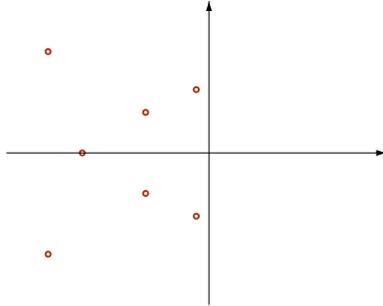


Рис. 4.3. Срыв устойчивости в динамической системе (4.7) происходит, когда собственные значения A «переходят» через мнимую ось.

то собственного значения, так что для определения этого момента нужно найти решение выписанного уравнения. Но такой ответ неправилен, так как для потери устойчивости необходимо не точное равенство нулю действительных частей некоторых собственных чисел матрицы, а переход их через нуль в область положительного знака. Без этого перехода через мнимую ось и «ещё чуть-чуть дальше» система останется устойчивой, сколь бы близко мы не придвинули собственные значения к мнимой оси или даже достигли бы её. Здесь важен именно переход «через и за» критическое значение, в отсутствие которого качественное изменение в поведении системы не совершится, и этот феномен совершенно не ухватывается понятиями ϵ -решения из §4.2в или ϵ -спектра из работ [11, 47, 59].

Рассмотренная ситуация, в действительности, весьма типична для динамических систем, где условием совершения многих типов структурных перестроек и изменений установившихся режимов работы систем — так называемых *бифуркаций* — является переход некоторого параметра через определённое *бифуркационное значение*. К примеру, при переходе через мнимую ось пары комплексных собственных чисел матрицы линеаризованной системы происходит бифуркация Андронова-Хопфа (называемая также «бифуркацией рождения цикла», см. [36]). И здесь принципиален именно переход через некоторый порог, а не близость к нему, на которую делается упор в понятиях ϵ -решения и ϵ -спектра.

Нетрудно понять, что такое «переход через нуль» для непрерывной функции одного переменного $f : \mathbb{R} \rightarrow \mathbb{R}$. Но в многомерной ситуа-

ции мы сталкиваемся с методическими трудностями, возникающими из необходимости иметь для нестрогого понятия «прохождение функции через нуль» чисто математическое определение. Из требования вычислительной корректности следует, что в любой окрестности такого решения каждая из компонент $F_i(x)$ вектор-функции $F(x)$ должна принимать как положительные, так и отрицательные значения. Но как именно? Какими должны (или могут) быть значения компонент $F_j(x)$, $j \neq i$, если $F_i(x) > 0$ или $F_i(x) < 0$?

В разрешении этого затруднения нам на помощь приходят нелинейный анализ и алгебраическая топология. В следующем параграфе мы приведём краткий набросок возможного решения этого вопроса.

4.3 Векторные поля и их вращение

4.3а Векторные поля

Если M — некоторое множество в \mathbb{R}^n и задано отображение

$$\Phi : M \rightarrow \mathbb{R}^n,$$

то часто удобно представлять значение $\Phi(x)$ как вектор, торчащий из точки $x \in M$. При этом говорят, что на M задано *векторное поле* Φ . Любопытно, что это понятие было введено около 1830 года М. Фарадеем, затем соответствующий язык проник в математическую физику, теорию дифференциальных уравнений и теорию динамических систем (см., к примеру, [8, 50]), и в настоящее время широко используется в современном естествознании. Мы воспользуемся соответствующими понятиями и результатами для наших целей анализа решений систем уравнений, численных методов и коррекции постановки задачи.

Векторное поле является *непрерывным*, если непрерывно отображение $\Phi(x)$. Например, на Рис. 4.4 изображены векторные поля

$$\Phi(x) = \Phi(x_1, x_2) = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \text{и} \quad \Psi(x) = \Psi(x_1, x_2) = \begin{pmatrix} x_1 \\ -x_2 \end{pmatrix}, \quad (4.8)$$

которые непрерывны и даже дифференцируемы.

Определение 4.3.1 Пусть задано векторное поле $\Phi : \mathbb{R}^n \supseteq M \rightarrow \mathbb{R}^n$. Точки $x \in M$, в которых поле обращается в нуль, т. е. $\Phi(x) = 0$, называются нулями поля или же его особыми точками.

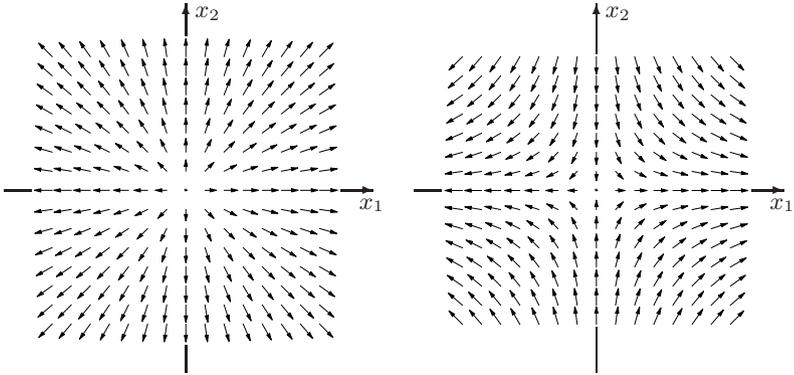


Рис. 4.4. Векторные поля $\Phi(x)$ и $\Psi(x)$, задаваемые формулами (4.8).

Связь векторных полей и их особых точек с основным предметом этой главы очевидна: особая точка поля $\Phi : M \rightarrow \mathbb{R}^n$ — это решение системы n уравнений

$$\left\{ \begin{array}{l} \Phi_1(x_1, x_2, \dots, x_n) = 0, \\ \Phi_2(x_1, x_2, \dots, x_n) = 0, \\ \quad \vdots \quad \ddots \quad \quad \quad \vdots \\ \Phi_n(x_1, x_2, \dots, x_n) = 0, \end{array} \right.$$

лежащее в M . Будем говорить, что векторное поле Φ вырождено, если у него есть особые точки. Иначе Φ называется невырожденным. К примеру, векторные поля Рис. 4.4 вырождены на всём \mathbb{R}^2 и имеют единственными особыми точками начало координат.

Определение 4.3.2 Пусть $\Phi(x)$ и $\Psi(x)$ — векторные поля на множестве $M \subseteq \mathbb{R}^n$. Непрерывная функция

$$\Delta(\lambda, x) : \mathbb{R} \times M \rightarrow \mathbb{R}^n$$

от параметра $\lambda \in [0, 1]$ и вектора $x \in \mathbb{R}^n$, такая что $\Phi(x) = \Delta(0, x)$ и $\Psi(x) = \Delta(1, x)$, называется деформацией векторного поля $\Phi(x)$ в векторное поле $\Psi(x)$.

Достаточно прозрачна связь деформаций с возмущениями векторного поля, т. е. отображения Φ . Но в качестве инструмента исследования решений систем уравнений и особых точек векторных полей нам нужны деформации, которые не искажают свойство поля быть невырожденным.

Определение 4.3.3 Деформацию $\Delta(\lambda, x)$ назовём невырожденной, если $\Delta(\lambda, x) \neq 0$ для всех $\lambda \in [0, 1]$ и $x \in M$.

Ясно, что невырожденные деформации могут преобразовывать друг в друга (соединять) только невырожденные векторные поля. Примерами невырожденных деформаций векторных полей, заданных на всём \mathbb{R}^n , являются растяжение, поворот относительно некоторой точки, параллельный перенос.

Определение 4.3.4 Если векторные поля можно соединить невырожденной деформацией, то они называются гомотопными.

В частности, любая достаточно малая деформация невырожденного векторного поля приводит к гомотопному полю.

Нетрудно понять, что отношение гомотопии векторных полей рефлексивно, симметрично и транзитивно, будучи поэтому *отношением эквивалентности*. Как следствие, непрерывные векторные поля, невырожденные на фиксированном множестве $M \subseteq \mathbb{R}^n$, распадается на классы гомотопных между собой полей.

4.3б Вращение векторных полей

Пусть D — ограниченная область в \mathbb{R}^n с границей ∂D . Через $\text{cl } D$ мы обозначим её топологическое замыкание. Оказывается, каждому невырожденному на ∂D векторному полю Φ можно сопоставить целочисленную характеристику — *вращение векторного поля Φ на ∂D* , — обозначаемую $\gamma(\Phi, D)$ и удовлетворяющую следующим условиям:

- (А) Гомотопные на ∂D векторные поля имеют одинаковое вращение.
- (В) Пусть $D_i, i = 1, 2, \dots$, — непересекающиеся области, лежащие в D (их может быть бесконечно много). Если непрерывное векторное поле Φ невырождено на теоретико-множественной разности

$$\text{cl } D \setminus \left(\bigcup_i D_i \right),$$

то вращения $\gamma(\Phi, D_i)$ отличны от нуля лишь для конечного набора D_i и

$$\gamma(\Phi, D) = \gamma(\Phi, D_1) + \gamma(\Phi, D_2) + \dots .$$

(С) Если $\Phi(x) = x - a$ для некоторой точки $a \in D$, то вращение Φ на ∂D равно $(+1)$, т. е.

$$\gamma(\Phi, D) = 1.$$

Нетрудно понять, что определенная так величина вращения поля устойчива к малым шевелениям как области (это следует из (В)), так и векторного поля (это вытекает из (А)).

Условиями (А)–(В)–(С) вращение векторного поля определяется однозначно, но можно показать [39], что это определение равносильно следующему конструктивному. Зафиксируем некоторую параметризацию поверхности ∂D , т. е. задание её в виде

$$\begin{aligned} x_1 &= x_1(u_1, u_2, \dots, u_{n-1}), \\ x_2 &= x_2(u_1, u_2, \dots, u_{n-1}), \\ &\vdots \quad \quad \quad \ddots \quad \quad \quad \vdots \\ x_n &= x_n(u_1, u_2, \dots, u_{n-1}), \end{aligned}$$

где u_1, u_2, \dots, u_{n-1} — параметры, $x_i(u_1, u_2, \dots, u_{n-1})$, $i = 1, 2, \dots, n$, — функции, определяющие одноименные координаты точки $x = (x_1, x_2, \dots, x_n) \in \partial D$. Тогда вращение поля $\Phi(x)$ на границе ∂D области D равно значению поверхностного интеграла

$$\frac{1}{S_n} \int_{\partial D} \frac{1}{\|\Phi(x)\|^n} \cdot \det \begin{pmatrix} \Phi_1(x) & \frac{\partial \Phi_1(x)}{\partial u_1} & \dots & \frac{\partial \Phi_1(x)}{\partial u_{n-1}} \\ \Phi_2(x) & \frac{\partial \Phi_2(x)}{\partial u_1} & \dots & \frac{\partial \Phi_2(x)}{\partial u_{n-1}} \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_n(x) & \frac{\partial \Phi_n(x)}{\partial u_1} & \dots & \frac{\partial \Phi_n(x)}{\partial u_{n-1}} \end{pmatrix} du_1 du_2 \dots du_n, \quad (4.9)$$

где S_n — площадь поверхности единичной сферы в \mathbb{R}^n . Этот интеграл обычно называют *интегралом Кронекера*.

В двумерном случае вращение векторного поля имеет простую геометрическую интерпретацию: это количество полных оборотов вектора

поля, совершаемое при движении точки аргумента в положительном направлении по рассматриваемой границе области [50, 54, 55, 57]. В многомерном случае такой наглядности уже нет, но величина вращения векторного поля Φ всё равно может быть истолкована как «число раз, которое отображение $\Phi : \partial D \rightarrow \Phi(\partial D)$ накрывает образ $\Phi(\partial D)$ ».

Рассмотрим примеры. На любой окружности с центром в нуле поле, изображенное на левой половине Рис. 4.4, имеет вращение $+1$, а поле на правой половине Рис. 4.4 — вращение -1 . Векторные поля Рис. 4.5, которые задаются формулами

$$\begin{cases} x_1 = r \cos(N\psi), \\ x_2 = r \sin(N\psi), \end{cases}$$

где $r = \sqrt{x_1^2 + x_2^2}$ — длина радиус-вектора точки $x = (x_1, x_2)$, ψ — его угол с положительным лучом оси абсцисс, при $N = 2$ и $N = 3$ имеют вращения $+2$ и $+3$ на окружностях с центром в нуле.

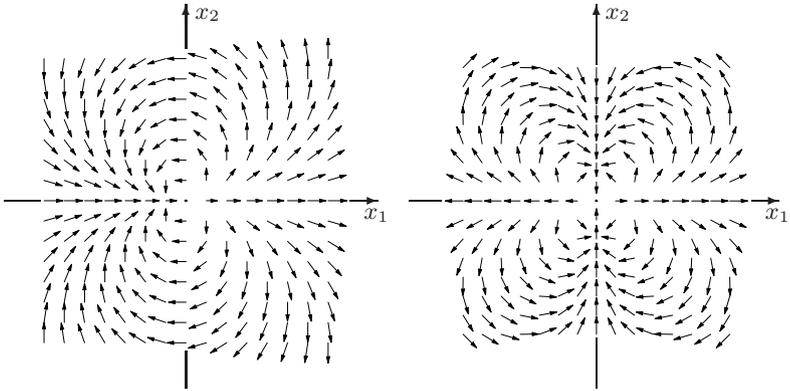


Рис. 4.5. Векторные поля, имеющие вращения $+2$ (левый чертёж) и $+3$ (правый чертёж) на любой окружности с центром в нуле.

С вращением векторного поля тесно связана другая известная глобальная характеристика отображений — *топологическая степень* [54, 55, 56, 57, 27, 39]. Именно, вращение поля Φ на границе области D есть

топологическая степень такого отображения ϕ границы ∂D в единичную сферу пространства \mathbb{R}^n , что

$$\phi(x) = \|\Phi(x)\|^{-1}\Phi(x).$$

Зачем нам понадобилось понятие вращения векторного поля? Мы собираемся использовать его для характеристики «прохождения через нуль» многомерной функции, и теоретической основой этого шага служат следующие результаты:

Предложение 4.3.1 [54, 55, 57] *Если векторное поле Φ невырождено на замыкании ограниченной области D , то вращение $\gamma(\Phi, D) = 0$.*

Теорема 4.3.1 (теорема Кронекера) [54, 55, 57] *Пусть векторное поле Φ невырождено на границе ограниченной области D и непрерывно на её замыкании. Если $\gamma(\Phi, D) \neq 0$, то поле Φ имеет в D по крайней мере одну особую точку.*

Теорема Кронекера обладает очень большой общностью и часто применяется не напрямую, а служит основой для более конкретных достаточных условий существования нулей поля или решений систем уравнений. Например, доказательство теоремы Миранды (см. §4.4б) сводится, фактически, к демонстрации того, что на границе области вращение векторного поля, соответствующего исследуемому отображению, равно ± 1 .

4.3в Индексы особых точек

Станем говорить, что особая точка является *изолированной*, если в некоторой её окрестности нет других особых точек рассматриваемого векторного поля. Таким образом, вращение поля одинаково на сферах достаточно малых радиусов с центром в изолированной особой точке \tilde{x} . Это общее вращение называют *индексом* особой точки \tilde{x} поля Φ или *индексом нуля \tilde{x}* поля Φ , и обозначают $\text{ind}(\tilde{x}, \Phi)$.

Итак, оказывается, что особые точки векторных полей (и решения систем уравнений) могут быть существенно разными, отличаясь друг от друга своим индексом, и различных типов особых точек существует столько же, сколько и целых чисел, т. е. счётное множество. Какими являются наиболее часто встречающиеся особые точки и, соответственно, решения систем уравнений? Ответ на этот вопрос даётся следующими двумя результатами:

Предложение 4.3.2 [54, 55, 57] *Если A — невырожденное линейное преобразование пространства \mathbb{R}^n , то его единственная особая точка — нуль — имеет индекс $\text{ind}(0, A) = \text{sgn det } A$.*

Определение 4.3.5 *Точка области определения отображения дифференцируемого отображения $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ называется критической, если в ней якобиан F' является особенной матрицей. Иначе говорят, что эта точка — регулярная.*

Предложение 4.3.3 [54, 55, 57] *Если \tilde{x} — регулярная особая точка дифференцируемого векторного поля Φ , то $\text{ind}(\tilde{x}, \Phi) = \text{sgn det } \Phi'(\tilde{x})$.*

Таким образом, регулярные (не критические) особые точки векторных полей имеют индекс ± 1 , а в прочих случаях значение индекса может быть весьма произвольным.

Например, индексы расположенного в начале координат нуля векторных полей, которые изображены на Рис. 4.4, равны $+1$ и (-1) , причём поля эти всюду дифференцируемы. Индексы нуля полей Рис. 4.5 равны $+2$ и $+3$, и в начале координат эти поля не дифференцируемы. Векторное поле на прямой, задаваемое рассмотренным в §4.26 квадратичным отображением $x \mapsto x^2 + px + q$ при $p^2 = 4q$ имеет особую точку $x = -p/2$ нулевого индекса.

4.3г Устойчивость особых точек

Определение 4.3.6 *Особая точка z поля Φ называется устойчивой, если для любого $\tau > 0$ можно найти такое $\eta > 0$, что всякое поле, отличающееся от Φ меньше чем на η , имеет особую точку, удаленную от z менее, чем на τ . Иначе особая точка z называется неустойчивой.*

Легко понять, что в связи с задачей решения систем уравнений нас интересуют именно устойчивые особые точки, поскольку задача поиска только таких точек является вычислительно-корректной.

Вторым основным результатом, ради которого мы затевали обзор теории вращения векторных полей, является следующее

Предложение 4.3.4 [55] *Изолированная особая точка непрерывного векторного поля устойчива тогда и только тогда, когда её индекс отличен от нуля.*

Например, неустойчивое решение квадратного уравнения (4.5)–(4.6) имеет индекс 0, а у векторных полей, изображённых на рисунках 4.4 и 4.5, начало координат является устойчивой особой точкой.

Интересно отметить, что отличие линейных уравнений от нелинейных, как следует из всего сказанного, проявляется не только в форме и структуре, но и в более глубоких вещах: 1) в линейных задачах индекс решения, как правило, равен ± 1 , а в нелинейных может быть как нулевым, так и отличным от ± 1 , и, как следствие, 2) в типичных линейных задачах изолированное решение устойчиво, а в нелинейных может быть неустойчивым.

Отметим отдельно, что результат об устойчивости особой точки ненулевого индекса ничего не говорит о количестве особых точек, близких к возмущаемой особой точке. В действительности, путем шевеления одной устойчивой особой точки можно получить сразу *несколько* особых точек, и это легко видеть на примере полей Рис. 4.5. Любая сколь угодно малая постоянная добавка к полю, изображённому на левом чертеже Рис. 4.5, приводит к распадению нулевой особой точки индекса 2 на две особые точки индекса 1. Аналогично, любая сколь угодно малая постоянная добавка к полю, изображённому на правом чертеже Рис. 4.5, приводит к распадению нулевой особой точки на три особые точки индекса 1. Таким образом, свойство единственности решения неустойчиво и требовать его наличия нужно со специальными оговорками.

Если в области D находится конечное число особых точек, то сумму их индексов называют *алгебраическим числом* особых точек.

Предложение 4.3.5 Пусть непрерывное векторное поле Φ имеет в D конечное число особых точек x_1, x_2, \dots, x_s и невырождено на границе ∂D . Тогда

$$\gamma(\Phi, D) = \text{ind}(x_1, \Phi) + \text{ind}(x_2, \Phi) + \dots + \text{ind}(x_s, \Phi).$$

Алгебраическое число особых точек устойчиво к малым возмущениям области и векторного поля, так как охватывает совокупную сумму индексов вне зависимости от рождения и уничтожения отдельных точек.

Наконец, сделаем ещё одно важное замечание. Нередко на практике для решения систем нелинейных уравнений исходную задачу переформулируют как оптимизационную, пользуясь, например, тем, что

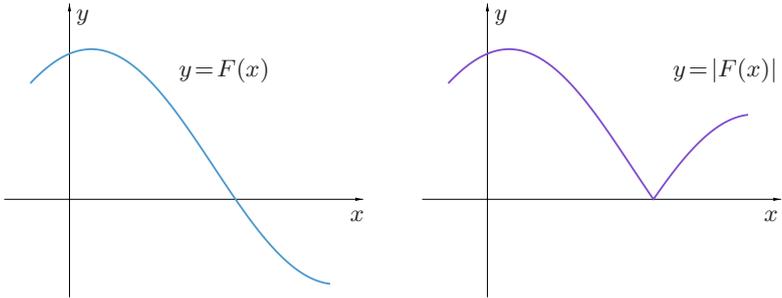


Рис. 4.6. Устойчивый нуль функции превращается в неустойчивый после взятия нормы функции.

справедливы следующие математические эквивалентности:

$$F(x) = 0 \quad \Leftrightarrow \quad \min_x \|F(x)\| = 0$$

и

$$F(x) = 0 \quad \Leftrightarrow \quad \min_x \|F(x)\|^2 = 0.$$

Далее имеющимися стандартными пакетами программ ищется решение задачи минимизации нормы $\|F(x)\|$ (или $\|F(x)\|^2$, чтобы обеспечить гладкость целевой функции) и результат сравнивается с нулём. С учётом наших знаний о задаче решения систем уравнений хорошо видна вычислительная неэквивалентность такого приведения: устойчивая особая точка *всегда* превращается при подобной трансформации в неустойчивое решение редуцированной задачи! Именно, любая сколь угодно малая добавка к $|F(x)|$ может приподнять график функции $y = |F(x)|$ над ось абсцисс (плоскостью нулевого уровня в общем случае), так что нуль функции исчезнет.

4.3д Вычислительно-корректная постановка задачи

Теперь все готово для вычислительно-корректной переформулировки задачи решения уравнений и систем уравнений. Она должна выглядеть следующим образом:

Для заданного $\varepsilon > 0$ и системы уравнений

$$F(x) = 0$$

найти на данном множестве $D \subseteq \mathbb{R}^n$

- 1) гарантированные двусторонние границы всех решений ненулевого индекса,
- 2) множество ε -решений.

(4.10)

Мы не требуем единственности решения в выдаваемых брусах, так как свойство решения быть единственным не является, как мы могли видеть, устойчивым к малым возмущениям задачи.

4.4 Классические методы решения уравнений

Пример 4.4.1 Рабочие имеют кусок кровельного материала шириной $l = 3.3$ метра и хотят покрыть им пролёт шириной $h = 3$ метра, сделав крышу круглой, в форме дуги окружности. Для того, чтобы придать правильную форму балкам, поддерживающим кровлю, нужно знать, какой именно радиус закругления крыши при этом получится (см. Рис. 4.7).

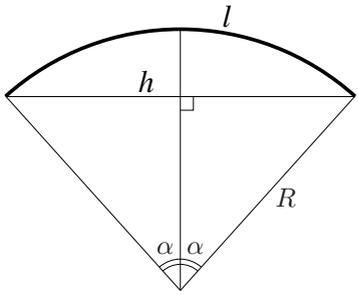


Рис. 4.7. Проектирование круглой крыши.

Обозначим искомый радиус закругления крыши через R . Если 2α — величина дуги (в радианах), соответствующей крыше, то

$$\frac{l}{2\alpha} = R.$$

С другой стороны, из рассмотрения прямоугольного треугольника с катетом $h/2$ и гипотенузой R получаем

$$R \sin \alpha = h/2.$$

Исключая из этих двух соотношений R , получим уравнение относительно одной неизвестной α :

$$l \sin \alpha = \alpha h. \quad (4.11)$$

Его решение не может быть выражено в явном виде, и потому далее мы обсудим возможности его численного решения. ■

Уравнение (4.11) является простейшим нелинейным *трансцендентным* уравнением. Так называют уравнения и системы уравнений, не являющиеся алгебраическими, т. е. такие, в которых в обеих частях уравнений стоят алгебраические выражения относительно неизвестных переменных.

4.4a Предварительная локализация решений

Обычно первым этапом численного решения уравнений и систем уравнений является предварительная локализация, т. е. уточнение местонахождения, искомых решений. Это вызвано тем, что большинство численных методов для поиска решений имеют локальный характер, т. е. сходятся к этим решениям лишь из достаточно близких начальных приближений.

Для локализации решений могут применяться как численные, так и аналитические методы, а также их смесь — гибридные методы, которые (следуя Д. Кнуту) можно назвать *получисленными* или *полуаналитическими*.

Особенно много аналитических результатов существует о локализации решений алгебраических уравнений (корней полиномов), что, конечно, имеет причину в очень специальном виде этих уравнений, допускающем исследование с помощью выкладок и т. п.

Теорема 4.4.1 Пусть для алгебраического уравнения вида

$$a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 = 0$$

обозначено

$$\alpha = \max\{a_0, \dots, a_{n-1}\}, \quad \beta = \max\{a_1, \dots, a_n\}.$$

Тогда все решения этого уравнения принадлежат кольцу в комплексной плоскости, определяемому условием

$$\frac{1}{1 + \beta/|a_0|} \leq |x| \leq 1 + \frac{\alpha}{|a_n|}.$$

Полезно правило знаков Декарта, утверждающее, что число положительных корней полинома с вещественными коэффициентами равно числу перемен знаков в ряду его коэффициентов или на чётное число меньше этого числа. При этом корни считаются с учётом кратности, а нулевые коэффициенты при подсчёте числа перемен знаков не учитываются. Если, к примеру, заранее известно, что все корни данного полинома вещественны, то правило знаков Декарта даёт точное число корней. Рассматривая полином с переменной $(-x)$ можно с помощью этого же результата найти число отрицательных корней исходного полинома.

4.46 Метод дихотомии

Этот метод часто называют также *методом бисекции* или *методом половинного деления*. Он заключается в последовательном делении пополам интервала локализации корня уравнения, на концах которого функция принимает значения разных знаков. Теоретической основой метода дихотомии является следующий факт, хорошо известный в математическом анализе:

Теорема 4.4.2 (теорема Больцано-Коши) Если функция $f : \mathbb{R} \rightarrow \mathbb{R}$ непрерывна на интервале $X \subset \mathbb{R}$ и на его концах принимает значения разных знаков, то внутри интервала X существует нуль функции f , т. е. точка $\tilde{x} \in X$, в которой $f(\tilde{x}) = 0$.

Часто её называют просто «теоремой Больцано» (см., к примеру, [38]), так как именно Б. Больцано первым обнаружил это замечательное свойство непрерывных функций.

Очевидно, что из двух половин интервала, на котором функция меняет знак, хотя бы на одной эта переменная знака обязана сохраняться. Её мы и оставляем в результате очередной итерации метода дихотомии.

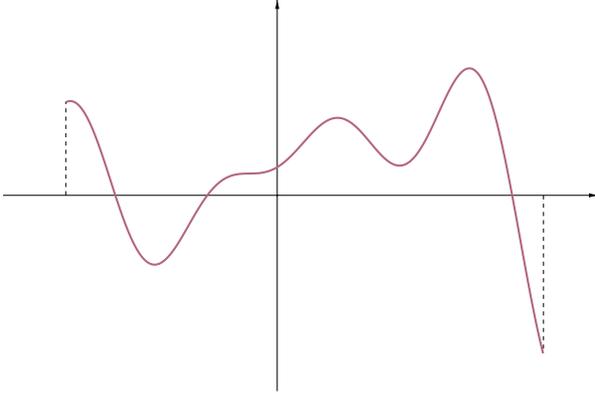


Рис. 4.8. Иллюстрация метода дихотомии (деления пополам)

На вход алгоритму подаются функция f , принимающая на концах интервала $[a, b]$ значения разных знаков, и точность ϵ , с которой необходимо локализовать решение уравнения $f(x) = 0$. На выходе получаем интервал $[\underline{x}, \bar{x}]$ шириной не более ϵ , содержащий решение уравнения.

Недостаток этого простейшего варианта метода дихотомии — возможность потери решений для функций, аналогичных изображенной на Рис. 4.8. На левой половине исходного интервала функция знака не меняет, но там находятся два нуля функции. Чтобы убедиться в единственности решения или в его отсутствии, можно привлекать дополнительную информацию об уравнении, к примеру, о производной фигурирующей в нём функции. В общем случае потери нулей можно избежать, если не отбрасывать подынтервалы, на которых доказательно не установлено отсутствие решений. Последовательная реализация этой идеи приводит к «методу ветвлений и отсечений», который подробно рассматривается далее в §4.8.

Многомерное обобщение теоремы Больцано-Коши было опубликовано более чем столетием позже в заметке [44]:

Теорема 4.4.3 (теорема Миранды) Пусть $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $f(x) = (f_1(x), f_2(x), \dots, f_n(x))^T$ — функция, непрерывная на бруске $X \subset \mathbb{R}^n$ со сторо-

Таблица 4.1. Метод дихотомии решения уравнений

$\underline{x} \leftarrow a; \quad \bar{x} \leftarrow b;$ DO WHILE $(\bar{x} - \underline{x} > \epsilon)$ $y \leftarrow \frac{1}{2}(\underline{x} + \bar{x});$ IF $(f(\underline{x}) < 0$ и $f(y) > 0)$ или $(f(\underline{x}) > 0$ и $f(y) < 0)$ $\bar{x} \leftarrow y$ ELSE $\underline{x} \leftarrow y$ END IF END DO
--

нами, параллельными координатным осям, и для любого $i = 1, 2, \dots, n$ имеет место либо

$$f_i(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \underline{\mathbf{X}}_i, \mathbf{X}_{i+1}, \dots, \mathbf{X}_n) \leq 0$$

$$\text{и } f_i(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \bar{\mathbf{X}}_i, \mathbf{X}_{i+1}, \dots, \mathbf{X}_n) \geq 0,$$

либо

$$f_i(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \underline{\mathbf{X}}_i, \mathbf{X}_{i+1}, \dots, \mathbf{X}_n) \geq 0$$

$$\text{и } f_i(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \bar{\mathbf{X}}_i, \mathbf{X}_{i+1}, \dots, \mathbf{X}_n) \leq 0,$$

т. е. области значений каждой компоненты функции $f(x)$ на соответствующих противоположных гранях бруса \mathbf{X} имеют разные знаки. Тогда на брус \mathbf{X} существует нуль функции f , т. е. точка $x^* \in \mathbf{X}$, в которой $f(x^*) = 0$.

Характерной особенностью теоремы Миранды является специальная форма множества, на котором утверждается существование нуля функции: оно должно быть брусом со сторонами, параллельными координатным осям, т. е. интервальным вектором. Для полноценного применения теоремы Миранды нужно уметь находить или как-то оценивать области значений функций на таких брусах.

Удобное средство для решения этой задачи предоставляют методы интервального анализа. Задача об определении области значений функции на брусках из области её определения эквивалентна задаче оптимизации, но в интервальном анализе она принимает специфическую форму задачи о вычислении так называемого *интервального расширения функции* (см. §1.5).

4.4в Метод простой итерации

Методом простой итерации обычно называют стационарный одношаговый итерационный процесс, который организуется после того, как исходное уравнение каким-либо способом приведено к равносильному рекуррентному виду $x = \Phi(x)$. Далее, после выбора некоторого начального приближения $x^{(0)}$, запускается итерационный процесс

$$x^{(k+1)} \leftarrow \Phi(x^{(k)}), \quad k = 0, 1, 2, \dots$$

При благоприятных обстоятельствах последовательность $\{x^{(k)}\}$ сходится, и её пределом является решение исходного уравнения. Но в общем случае и характер сходимости, и вообще её наличие существенно зависят как от отображения Φ , так и от начального приближения к решению.

Пример 4.4.2 Уравнение (4.11) из Примера 4.4 нетрудно привести к рекуррентному виду

$$\alpha = \frac{l}{h} \sin \alpha,$$

где $l = 3.3$ и $h = 3$. Далее, взяв в качестве начального приближения, например, $\alpha^{(0)} = 1$, через 50 итераций

$$\alpha^{(k+1)} \leftarrow \frac{l}{h} \sin \alpha^{(k)}, \quad k = 0, 1, 2, \dots, \quad (4.12)$$

получаем пять верных знаков точного решения $\alpha^* = 0.748986642697\dots$ (читатель легко может самостоятельно проверить все числовые данные этого примера с помощью любой системы компьютерной математики).

Итерационный процесс (4.12) сходится к решению α^* не из любого начального приближения. Если $\alpha^{(0)} = \pi l$, $l \in \mathbb{Z}$, то выполнение итераций (4.12) с идеальной точностью даёт $\alpha^{(k)} = 0$, $k = 1, 2, \dots$. Если же $\alpha^{(0)}$ таково, что синус от него отрицателен, то итерации (4.12) сходятся к

решению $(-\alpha^*)$ уравнения (4.11). И нулевое, и отрицательное решения очевидно не имеют содержательного смысла.

С другой стороны, переписывание исходного уравнения (4.12) в другом рекуррентном виде —

$$\alpha = \frac{1}{l} \arcsin(\alpha h)$$

— приводит к тому, что характер сходимости метода простой итерации совершенно меняется. Из любого начального приближения, меньшего по модулю чем примерно 0.226965, итерации

$$\alpha^{(k+1)} \leftarrow \frac{1}{l} \arcsin(\alpha^{(k)} h), \quad k = 0, 1, 2, \dots,$$

сходятся лишь к нулевому решению. Большие по модулю начальные приближения быстро выводят за границы области определения вещественного арксинуса, переводя итерации в комплексную плоскость, где они снова сходятся к нулевому решению. Таким образом, искомого решения α^* мы при этом никак не получаем. ■

Рассмотренный пример хорошо иллюстрирует различный характер неподвижных точек отображений и мотивирует следующие определения.

Неподвижная точка x^* функции $\Phi(x)$ называется *притягивающей*, если существует такая окрестность Ω точки x^* , что итерационный процесс $x^{(k+1)} \leftarrow \Phi(x^{(k)})$ сходится к x^* из любого начального приближения $x^{(0)} \in \Omega$.

Неподвижная точка x^* функции $\Phi(x)$ называется *отталкивающей*, если существует такая окрестность Ω точки x^* , что итерационный процесс $x^{(k+1)} \leftarrow \Phi(x^{(k)})$ не сходится к x^* при любом начальном приближении $x^{(0)} \in \Omega$.

Ясно, что простые итерации $x^{(k+1)} \leftarrow \Phi(x^{(k)})$ непригодны для нахождения отталкивающих неподвижных точек. Здесь возникает интересный вопрос о том, какими преобразованиями уравнений и систем уравнений отталкивающие точки можно сделать притягивающими.

Наиболее часто существование неподвижных точек можно гарантировать у отображений, которые удовлетворяют тем или иным дополнительным условиям, и самыми популярными из них являются так называемые условия сжимаемости (сжатия) образа.

Напомним, что отображение $g : X \rightarrow X$ метрического пространства X с расстоянием $\text{dist} : X \rightarrow \mathbb{R}_+$ называется *сжимающим* (или просто *сжатием*), если существует такая положительная постоянная $\alpha < 1$, что для любой пары элементов $x, y \in X$ имеет место неравенство

$$\text{dist}(g(x), g(y)) \leq \alpha \cdot \text{dist}(x, y).$$

Теорема 4.4.4 (теорема Банаха о неподвижной точке). *Сжимающее отображение $g : X \rightarrow X$ полного метрического пространства X в себя имеет единственную неподвижную точку. Она может быть найдена методом последовательных приближений*

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

при любом начальном приближении $x^{(0)} \in X$.

Доказательство этого результата можно найти, к примеру, в [17]. Особенно ценен в теореме Банаха её конструктивный характер, позволяющий организовать численные методы для нахождения неподвижной точки.

Иногда бывает полезно работать с векторнозначным расстоянием — *мультиметрикой*, — которая вводится на \mathbb{R}^n как

$$\text{Dist}(x, y) := \begin{pmatrix} \text{dist}(x_1, y_1) \\ \vdots \\ \text{dist}(x_n, y_n) \end{pmatrix} \in \mathbb{R}_+^n. \quad (4.13)$$

Для мультиметрических пространств аналогом теоремы Банаха о неподвижной точке для сжимающих отображений является приводимая ниже теорема Шрёдера о неподвижной точке. Перед тем, как дать её точную формулировку, введём

Определение 4.4.1 *Отображение $g : X \rightarrow X$ мультиметрического пространства X с мультиметрикой $\text{Dist} : X \rightarrow \mathbb{R}_+^n$ называется P -сжимающим (или просто P -сжатием), если существует неотрицательная $n \times n$ -матрица P со спектральным радиусом $\rho(P) < 1$, такая что для всех $x, y \in X$ имеет место*

$$\text{Dist}(g(x), g(y)) \leq P \cdot \text{Dist}(x, y). \quad (4.14)$$

Следует отметить, что математики, к сожалению, не придерживаются здесь единой терминологии. Ряд авторов (см. [46]) за матрицей P из (4.14) закрепляют отдельное понятие «оператора Липшица (матрицы Липшица) отображения g », и в условиях Определения 4.4.1 говорят, что «оператор Липшица для g сжимающий».

Теорема 4.4.5 (теорема Шрёдера о неподвижной точке) *Пусть отображение $g : \mathbb{R}^n \supseteq X \rightarrow \mathbb{R}^n$ является P -сжимающим на замкнутом подмножестве X пространства \mathbb{R}^n с мультиметрикой Dist . Тогда для любого $x^{(0)}$ последовательность итераций*

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

сходится к единственной неподвижной точке x^ отображения g в X и имеет место оценка*

$$\text{Dist}(x^{(k)}, x^*) \leq (I - P)^{-1} P \cdot \text{Dist}(x^{(k)}, x^{(k-1)}).$$

Доказательство можно найти, например, в книгах [1, 18, 27, 46]

4.4г Метод Ньютона и его модификации

Предположим, что для уравнения $f(x) = 0$ с вещественнозначной функцией f известно некоторое приближение \tilde{x} к решению x^* . Если f — плавно меняющаяся (гладкая функция), то естественно приблизить её в окрестности точки \tilde{x} линейной функцией, т. е.

$$f(x) \approx f(\tilde{x}) + f'(\tilde{x})(x - \tilde{x}),$$

и далее для вычисления следующего приближения к x^* решать линейное уравнение

$$f(\tilde{x}) + f'(\tilde{x})(x - \tilde{x}) = 0.$$

Отсюда очередное приближение к решению

$$x = \tilde{x} - \frac{f(\tilde{x})}{f'(\tilde{x})}.$$

Итерационный процесс

$$x^{(k+1)} \leftarrow x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k = 0, 1, 2, \dots,$$

называют *методом Ньютона*. Он является одним из наиболее популярных и наиболее эффективных численных методов решения уравнений и имеет многочисленные обобщения, в том числе на многомерный случай, т. е. в применении к решению систем уравнений (см. 4.7в).

Пример 4.4.3 Рассмотрим уравнение $x^2 - a = 0$, решением которого является квадратный корень из числа a . Если $f(x) = x^2 - a$, то $f'(x) = 2x$, так что в методе Ньютона для нахождения решения рассматриваемого уравнения имеем

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})} = x^{(k)} - \frac{(x^{(k)})^2 - a}{2x^{(k)}} = \frac{x^{(k)}}{2} + \frac{a}{2x^{(k)}}.$$

Итерационный процесс для нахождения квадратного корня

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

известен ещё с античности и часто называется *методом Герона*. Для любого положительного начального приближения $x^{(0)}$ он порождает убывающую, начиная с $x^{(1)}$, последовательность, которая быстро сходится к арифметическому значению \sqrt{a} . ■

Метод Ньютона требует вычисления на каждом шаге производной от функции f , что может оказаться неприемлемым или труднодостижимым. Одна из очевидных модификаций метода Ньютона состоит в том, чтобы «заморозить» производную в некоторой точке и вести итерации по формуле

$$x^{(k+1)} \leftarrow x^{(k)} - \frac{f(x^{(k)})}{f'(\tilde{x})}, \quad k = 0, 1, 2, \dots,$$

где \tilde{x} — фиксированная точка, в которой берётся производная. Получаем стационарный итерационный процесс, который существенно проще в реализации, но он имеет качественно более медленную сходимость.

Для определения погрешности приближённого решения \tilde{x} и контроля точности вычислений можно применять формулу

$$|\tilde{x} - x^*| \leq \frac{|f(\tilde{x})|}{\min_{\xi \in [a, b]} |f'(\xi)|}, \quad (4.15)$$

которая следует из теоремы Лагранжа о среднем (формулы конечных приращений):

$$f(\tilde{x}) - f(x^*) = f'(\xi)(\tilde{x} - x^*),$$

где ξ — некоторая точка, заключённая между \tilde{x} и x^* , т. е. $\xi \in \square\{\tilde{x}, x^*\}$. Ясно, что тогда

$$|f(\tilde{x}) - f(x^*)| \geq \min_{\xi} |f'(\xi)| \cdot |\tilde{x} - x^*|,$$

и при $\min_{\xi \in [a, b]} |f'(\xi)| \neq 0$ получаем оценку (4.15). Отметим её очевидную аналогию с оценкой (3.130) для погрешности решения систем линейных уравнений.

4.4д Методы Чебышёва

Методы Чебышёва для решения уравнения $f(x) = 0$ основаны на разложении по формуле Тейлора функции f^{-1} , обратной к f . Они могут иметь произвольно высокий порядок точности, определяемый количеством членов разложения для f^{-1} , но практически обычно ограничиваются небольшими порядками.

Предположим, что вещественная функция f является гладкой и монотонной на интервале $[a, b]$, так что она взаимно однозначно отображает этот интервал в некоторый интервал $[\alpha, \beta]$. Как следствие, существует обратная к f функция $g = f^{-1} : [\alpha, \beta] \rightarrow [a, b]$, которая имеет ту же гладкость, что и функция f .

Итак, пусть известно некоторое приближение \tilde{x} к решению x^* уравнения $f(x) = 0$. Обозначив $y = f(\tilde{x})$, разложим обратную функцию g в точке \tilde{x} по формуле Тейлора с остаточным членом в форме Лагранжа:

$$\begin{aligned} g(0) &= g(y) + g'(y)(0 - y) + g''(y) \frac{(0 - y)^2}{2} + \dots + g^{(p)}(y) \frac{(0 - y)^p}{p!} \\ &\quad + g^{(p+1)}(\xi) \frac{(0 - y)^{p+1}}{(p+1)!} \\ &= g(y) + \sum_{k=1}^p (-1)^k g^{(k)}(y) \frac{y^k}{k!} + (-1)^{p+1} g^{(p+1)}(\xi) \frac{y^{p+1}}{(p+1)!}, \end{aligned}$$

где ξ — какая-то точка между 0 и y . Возвращаясь к переменной x ,

будем иметь

$$x^* = \tilde{x} + \sum_{k=1}^p (-1)^k g^{(k)}(f(\tilde{x})) \frac{(f(\tilde{x}))^k}{k!} + (-1)^{p+1} g^{(p+1)}(\xi) \frac{(f(\tilde{x}))^{p+1}}{(p+1)!}.$$

В качестве следующего приближения к решению мы можем взять, отбросив остаточный член, значение

$$\tilde{x} + \sum_{k=1}^p (-1)^k g^{(k)}(f(\tilde{x})) \frac{(f(\tilde{x}))^k}{k!}.$$

Подытоживая сказанное, определим итерации

$$x^{(k+1)} \leftarrow x^{(k)} + \sum_{k=1}^p (-1)^k g^{(k)}(f(\tilde{x})) \frac{(f(\tilde{x}))^k}{k!}, \quad k = 0, 1, 2, \dots,$$

которые и называются *методом Чебышёва* p -го порядка.

Как найти производные обратной функции g ?

Зная производные функции f в какой-то точке, мы можем найти и производные обратной функции g . В самом деле, последовательно дифференцируя тождество $x = g(f(x))$, получим

$$\begin{aligned} g'(f(x)) f'(x) &= 1, \\ g''(f(x)) (f'(x))^2 + g'(f(x)) f''(x) &= 0, \\ g'''(f(x)) (f'(x))^3 + g''(f(x)) \cdot 2f'(x)f''(x) \\ &\quad + g'(f(x)) f'(x) f'''(x) + g'(f(x)) f'''(x) = 0, \\ &\dots \qquad \qquad \qquad \dots \end{aligned}$$

или

$$\begin{aligned} g'(f(x)) f'(x) &= 1, \\ g''(f(x)) (f'(x))^2 + g'(f(x)) f''(x) &= 0, \\ g'''(f(x)) (f'(x))^3 + 3g''(f(x)) f'(x) f''(x) + g'(f(x)) f'''(x) &= 0, \\ &\dots \qquad \qquad \qquad \dots \end{aligned}$$

Относительно неизвестных значений производных $g'(f(x))$, $g''(f(x))$, $g'''(f(x))$ и т. д. эта система соотношений имеет треугольный вид, позволяющий найти их последовательно одну за другой:

$$\begin{aligned} g'(f(x)) &= \frac{1}{f'(x)}, \\ g''(f(x)) &= -\frac{g(f(x)) f''(x)}{(f'(x))^2} = -\frac{f''(x)}{(f'(x))^3}, \\ g'''(f(x)) &= -\frac{3g''(f(x)) f'(x) f''(x) + g'(f(x)) f'''(x)}{(f'(x))^3} \\ &= -3 \frac{(f''(x))^2}{(f'(x))^5} - \frac{f'''(x)}{(f'(x))^4} \end{aligned}$$

и так далее.

Для $p = 1$ расчётные формулы метода Чебышёва имеют вид

$$x^{(k+1)} \leftarrow x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k = 0, 1, 2, \dots,$$

что совпадает с методом Ньютона.

Для $p = 2$ расчётные формулы метода Чебышёва таковы

$$x^{(k+1)} \leftarrow x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})} - \frac{f''(x^{(k)}) (f(x^{(k)}))^2}{2(f'(x^{(k)}))^3}, \quad k = 0, 1, 2, \dots$$

Наиболее часто методом Чебышёва называют именно этот итерационный процесс, так как методы более высокого порядка из этого семейства на практике используются редко.

4.5 Классические методы решения систем уравнений

4.5а Метод простой итерации

Схема применения метода простой итерации для систем уравнений в принципе не отличается от случая одного уравнения. Исходная система уравнений $F(x) = 0$ должна быть каким-либо способом приведена

Итерационный процесс

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

называют *методом Ньютона*.

Метод Ньютона требует вычисления на каждом шаге матрицы производных функции F и решения систем линейных алгебраических уравнений с изменяющейся матрицей. Нередко подобные трудозатраты могут стать излишне обременительными. Если зафиксировать точку \tilde{x} , в которой вычисляется эта матрица производных, то получим упрощённый стационарный итерационный процесс

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

который часто называют *модифицированным методом Ньютона*. Здесь решение систем линейных уравнений с одинаковыми матрицами $F'(\tilde{x})$ можно осуществлять по упрощённым алгоритмам, к примеру, найдя один раз LU-разложение матрицы и далее используя его.

Один из наиболее часто используемых результатов о сходимости метода Ньютона — это

Теорема 4.5.1 (теорема Канторовича о методе Ньютона)

Пусть отображение F определено в открытой области D и имеет непрерывную вторую производную F'' в замыкании $\text{cl } D$. Пусть, кроме того, существует такой непрерывный линейный оператор $\Gamma_0 = (F'(x_0))^{-1}$, что $\|\Gamma_0 F(x_0)\| \leq \eta$ и $\|\Gamma_0 F''(x)\| < K$ для всех $x \in \text{cl } D$ и некоторых констант η и K . Если

$$h = K\eta \leq \frac{1}{2}$$

и

$$r \geq r_0 = \frac{1 - \sqrt{1 - 2h}}{h} \eta,$$

то уравнение $F(x) = 0$ имеет решение x^ , к которому сходится метод Ньютона, как исходный, так и модифицированный. При этом*

$$\|x^* - x_0\| \leq r_0.$$

Для исходного метода Ньютона сходимость описывается оценкой

$$\|x^* - x_k\| \leq \frac{\eta}{2^k h} (2h)^{2^k}, \quad k = 0, 1, 2, \dots,$$

а для модифицированного метода верна оценка

$$\|x^* - x_k\| \leq \frac{\eta}{h}(1 - \sqrt{1 - 2h})^{k+1}, \quad k = 0, 1, 2, \dots,$$

при условии $h < \frac{1}{2}$.

Доказательство и дальнейшие результаты на эту тему можно найти в книге [17].

4.6 Интервальные линейные системы уравнений

Предметом рассмотрения настоящего пункта являются интервальные системы линейных алгебраических уравнений (ИСЛАУ) вида

$$Ax = b, \quad (4.16)$$

где $A = (a_{ij})$ — это интервальная $m \times n$ -матрица и $b = (b_i)$ — интервальный m -вектор. Для интервальных уравнений решения и множества решений могут быть определены разнообразными способами (см. [37]), но ниже мы ограничимся так называемым *объединённым множеством решений* для (4.16), которое образовано всевозможными решениями x точечных систем $Ax = b$, когда матрица A и вектор b независимо пробегают \mathbf{A} и \mathbf{b} соответственно. Объединённое множество решений определяется строго как

$$\Xi(\mathbf{A}, \mathbf{b}) = \{x \in \mathbb{R}^n \mid (\exists A \in \mathbf{A})(\exists b \in \mathbf{b})(Ax = b)\}, \quad (4.17)$$

и ниже мы будем называть его просто *множеством решений* интервальной линейной системы (4.16), так как другие множества решений нами не исследуются. Точное описание множества решений может расти экспоненциально с размерностью вектора неизвестных n , а потому является практически невозможным уже при n , превосходящем несколько десятков. С другой стороны, в большинстве реальных постановок задач точное описание на самом деле и не нужно. На практике бывает вполне достаточно нахождения *оценки* для множества решений, т. е. приближенного описания, удовлетворяющего содержательно-смыслу рассматриваемой задачи.

Приведём полезный технический результат, который часто используется в связи с исследованием и оцениванием множества решений интервальных систем линейных алгебраических уравнений.

Теорема 4.6.1 (характеризация Бека) *Если $\mathbf{A} \in \mathbb{IR}^{m \times n}$, $\mathbf{b} \in \mathbb{IR}^m$, то*

$$\begin{aligned}\Xi(\mathbf{A}, \mathbf{b}) &= \{x \in \mathbb{R}^n \mid \mathbf{A}x \cap \mathbf{b} \neq \emptyset\} \\ &= \{x \in \mathbb{R}^n \mid 0 \in \mathbf{A}x - \mathbf{b}\}.\end{aligned}$$

Доказательство. Если $\tilde{x} \in \Xi(\mathbf{A}, \mathbf{b})$, то $\tilde{A}\tilde{x} = \tilde{b}$ для некоторых $\tilde{A} \in \mathbf{A}$, $\tilde{b} \in \mathbf{b}$. Следовательно, по крайней мере $\tilde{b} \in \mathbf{A}\tilde{x} \cap \mathbf{b}$, так что действительно $\mathbf{A}\tilde{x} \cap \mathbf{b} \neq \emptyset$.

Наоборот, если $\mathbf{A}\tilde{x} \cap \mathbf{b} \neq \emptyset$, то это пересечение $\mathbf{A}\tilde{x} \cap \mathbf{b}$ содержит вектор $\tilde{b} \in \mathbb{R}^m$, для которого должно иметь место равенство $\tilde{b} = \tilde{A}\tilde{x}$ с некоторой $\tilde{A} \in \mathbf{A}$. Итак, $\tilde{x} \in \Xi(\mathbf{A}, \mathbf{b})$.

Второе равенство следует из того, что $\mathbf{A}\tilde{x} \cap \mathbf{b} \neq \emptyset$ тогда и только тогда, когда $0 \in \mathbf{A}\tilde{x} - \mathbf{b}$. ■

Теорема 4.6.2 (характеризация Оеттли-Прагера) *Для объединённого множества решений ИСЛАУ имеет место*

$$x \in \Xi(\mathbf{A}, \mathbf{b}) \Leftrightarrow |(\text{mid } \mathbf{A})x - \text{mid } \mathbf{b}| \leq \text{rad } \mathbf{A} \cdot |x| + \text{rad } \mathbf{b}, \quad (4.18)$$

где неравенство между векторами понимается покомпонентным образом.

Доказательство. Для любых интервальных векторов-брусков \mathbf{p} и \mathbf{q} включение $\mathbf{p} \subseteq \mathbf{q}$ равносильно покомпонентному неравенству

$$|\text{mid } \mathbf{q} - \text{mid } \mathbf{p}| \leq \text{rad } \mathbf{q} - \text{rad } \mathbf{p}.$$

Следовательно, условие характеристики Бека, т. е. $0 \in \mathbf{A}\tilde{x} - \mathbf{b}$, может быть переписано в следующем виде:

$$|\text{mid } (\mathbf{A}\tilde{x} - \mathbf{b})| \leq \text{rad } (\mathbf{A}\tilde{x} - \mathbf{b}).$$

С учётом правил преобразования середины и радиуса получаем

$$\begin{aligned}\text{mid } (\mathbf{A}\tilde{x} - \mathbf{b}) &= (\text{mid } \mathbf{A})\tilde{x} - \text{mid } \mathbf{b}, \\ \text{rad } (\mathbf{A}\tilde{x} - \mathbf{b}) &= (\text{rad } \mathbf{A}) \cdot |\tilde{x}| + \text{rad } \mathbf{b},\end{aligned}$$

откуда вытекает требуемое. ■

4.7 Интервальные методы решения уравнений и систем уравнений

Интервальные методы позволяют придать конструктивный характер некоторым известным результатам математического анализа, которые раньше рассматривались как «чистые» теоремы существования. Самым первым из них является

Теорема 4.7.1 (теорема Брауэра о неподвижной точке) *Пусть D — выпуклое компактное множество в \mathbb{R}^n . Если непрерывное отображение $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ переводит D в себя, $g(D) \subseteq D$, то оно имеет на D неподвижную точку x^* , т. е. такую что $x^* = g(x^*)$.*

Если вместо произвольных выпуклых компактов ограничиться интервальными векторами-брусами в \mathbb{R}^n , а для оценивания области значений применять его внешнюю оценку в виде интервального расширения, то условия теоремы Брауэра могут быть конструктивно проверены на компьютере.

С учётом сказанного выше во введении к главе (стр. 435) о равносильности рекуррентного вида систем уравнений (4.4) канонической форме (4.1)–(4.2) чрезвычайно полезными для вычислительной математики оказываются результаты анализа, утверждающие существование неподвижных точек отображений. Теорема Брауэра является именно таким результатом.

4.7а Основы интервальной техники

Задача решения уравнений и систем уравнений является одной из классических задач вычислительной математики, для решения которой развито немало эффективных подходов — метод простой итерации, метод Ньютона, их модификации и т.п. Преимущества и недостатки этих классических методов мы обсудили выше в §§4.4–4.5 (см. также [5, 27, 31, 42]). Для дальнейшего нам важны два факта:

- Для уравнений, в которых фигурируют функции, не обладающие «хорошими» глобальными свойствами, все традиционные методы имеют *локальный характер*, т. е. обеспечивают отыскание решения, находящегося в некоторой (иногда достаточно малой) окрестности начального приближения. Задача нахождения *всех*

решений уравнения или системы уравнений, как правило, рассматривается лишь в специальных руководствах и методы её решения оказываются очень сложными.

- Гарантированные оценки погрешности найденного приближения к решению в традиционных методах дать весьма непросто.

Указание приближённого значения величины и его максимальной погрешности равносильно тому, что мы знаем левую и правую границы возможных значений этой величины, и поэтому можно переформулировать нашу задачу в следующем усиленном виде —

<p>Для каждого решения системы уравнений</p> $F(x) = 0$ <p>на данном множестве $D \subseteq \mathbb{R}^n$ найти гарантированные двусторонние границы</p>), (4.19)
---	-----------

— который будем называть *задачей доказательного глобального решения* системы уравнений. Эпитет «доказательный» означает здесь, что получаемый нами ответ к задаче — границы решений и т.п. — имеет статус математически строго доказанного утверждения о расположении решений при условии, что ЭВМ работает корректно (см. §1.8).

Задача (4.19) оказывается чрезвычайно сложной, а в классическом численном анализе почти полностью отсутствуют развитые методы для её решения. Из часто используемых подходов, имеющих ограниченный успех, следует упомянуть *аналитическое исследование, мультистарт, методы продолжения* [27].

Итак, пусть к решению предъявлена система уравнений (4.2)

$$F(x) = 0$$

на бруске $\mathbf{X} \subset \mathbb{R}^n$. Существование решения этой системы на \mathbf{X} можно переписать в виде равносильного условия

$$\text{ran}(F, \mathbf{X}) \ni 0,$$

и потому техника интервального оценивания множеств значений функций оказывается весьма полезной при решении рассматриваемой задачи. В частности, если нуль содержится во внутренней интервальной оценке множества значений $\text{ran}(F, \mathbf{X})$ отображения F , то на бруске

\mathbf{X} гарантированно находится решение системы (4.2). С другой стороны, если в нашем распоряжении имеется интервальное расширение \mathbf{F} функции F на \mathbf{X} , то $\mathbf{F}(\mathbf{X}) \supseteq \text{ran}(F, \mathbf{X})$. Поэтому если $0 \notin \mathbf{F}(\mathbf{X})$, то на \mathbf{X} нет решений рассматриваемой системы уравнений.

Далее, перепишем исходную систему (4.2) в равносильной рекуррентной форме

$$x = T(x) \quad (4.20)$$

с некоторым отображением $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$. Оно может быть взято, как примеру, в виде

$$T(x) = x - F(x)$$

либо

$$T(x) = x - AF(x),$$

с неособенной $n \times n$ -матрицей A , либо как-нибудь ещё. Пусть также $\mathbf{T} : \mathbb{IR}^n \rightarrow \mathbb{IR}^n$ — интервальное расширение отображения T . Ясно, что решения системы (4.20) могут лежать лишь в пересечении $\mathbf{X} \cap \mathbf{T}(\mathbf{X})$. Поэтому если

$$\mathbf{X} \cap \mathbf{T}(\mathbf{X}) = \emptyset,$$

то в \mathbf{X} нет решений системы уравнений (4.20). Коль скоро искомое решение содержится и в $\mathbf{T}(\mathbf{X})$, то для дальнейшего уточнения бруса, в котором может присутствовать решение, мы можем организовать итерации с пересечением

$$\mathbf{X}^{(0)} \leftarrow \mathbf{X}, \quad (4.21)$$

$$\mathbf{X}^{(k+1)} \leftarrow \mathbf{T}(\mathbf{X}^{(k)}) \cap \mathbf{X}^{(k)}, \quad k = 0, 1, 2, \dots \quad (4.22)$$

Следует особо отметить, что в получающихся при этом брусах наличие решения, вообще говоря, не гарантируется. Они являются лишь «подозрительными» на существование решения.

Но вот если для бруса \mathbf{X} выполнено

$$\mathbf{T}(\mathbf{X}) \subseteq \mathbf{X},$$

то по теореме Брауэра о неподвижной точке (стр. 472) в \mathbf{X} гарантированно находится решение системы (4.20). Для уточнения этого бруса мы снова можем воспользоваться итерациями (4.21)–(4.22). Таким образом, наихудшим, с точки зрения уточнения информации о решении системы, является случай

$$\mathbf{T}(\mathbf{X}) \supsetneq \mathbf{X}. \quad (4.23)$$

Приведённую выше последовательность действий по обнаружению решения системы уравнений и уточнению его границ мы будем называть далее кратко *тестом существования* (решения). Условимся также считать, что его результатом является брус пересечения ($\mathbf{X} \cap \mathbf{T}(\mathbf{X})$) либо предел последовательности (4.21)–(4.22). Если этот брус непуст, то он либо наверняка содержит решение системы уравнений, либо является подозрительным на наличие в нём решения. Если же результат теста существования пуст, то в исходном брусе решений системы уравнений нет.

В действительности, каждый из изложенных выше приёмов уточнения решения допускает далеко идущие модификации и улучшения. Например, это относится к итерациям вида (4.21)–(4.22), которые могут быть последовательно применены не к целым брусам $\mathbf{X}^{(k)}$, а к отдельным их компонентам в комбинации с различными способами приведения исходной системы к рекуррентному виду (4.20). На этом пути мы приходим к чрезвычайно эффективным алгоритмам, которые получили наименование *методов распространения ограничений* (см., к примеру, [30]).

Как простейший тест существования, так и его более продвинутые варианты без особых проблем реализуются на ЭВМ и работают тем лучше, чем более качественно вычисляются интервальные расширения функций F в (4.2) и T в (4.20) и чем меньше ширина бруса \mathbf{X} . Последнее связано с тем, что погрешность оценивания области значений функции посредством любого интервального расширения убывает с уменьшением размеров бруса, на котором производится это оценивание. (см. §1.5).

4.76 Одномерный интервальный метод Ньютона

В этом параграфе мы рассмотрим простейший случай одного уравнения с одним неизвестным.

Предположим, что $f : \mathbb{R} \supseteq \mathbf{x} \rightarrow \mathbb{R}$ — непрерывно дифференцируемая функция, имеющая нуль x^* на интервале \mathbf{x} , т.е. $f(x^*) = 0$. Тогда для любой точки $\tilde{x} \in \mathbf{x}$ из этого же интервала в силу теоремы Лагранжа о среднем значении

$$f(\tilde{x}) - f(x^*) = (\tilde{x} - x^*) \cdot f'(\xi),$$

где ξ — некоторая точка между \tilde{x} и x^* . Но так как $f(x^*) = 0$, то отсюда

следует

$$x^* = \tilde{x} - \frac{f(\tilde{x})}{f'(\xi)}.$$

Если $f'(x)$ является каким-либо интервальным расширением производной функции $f(x)$ на x , то $f'(\xi) \in f'(x)$ и

$$x^* \in \tilde{x} - \frac{f(\tilde{x})}{f'(x)}.$$

Интервальное выражение, фигурирующее в правой части этого включения, будет играть в дальнейшем важную роль и потому достойно выделения самостоятельным понятием.

Определение 4.7.1 Для заданной функции $f: \mathbb{R} \rightarrow \mathbb{R}$ отображение

$$\mathcal{N}: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R},$$

действующее по правилу

$$\mathcal{N}(x, \tilde{x}) := \tilde{x} - \frac{f(\tilde{x})}{f'(x)}$$

называется (одномерным) интервальным оператором Ньютона.

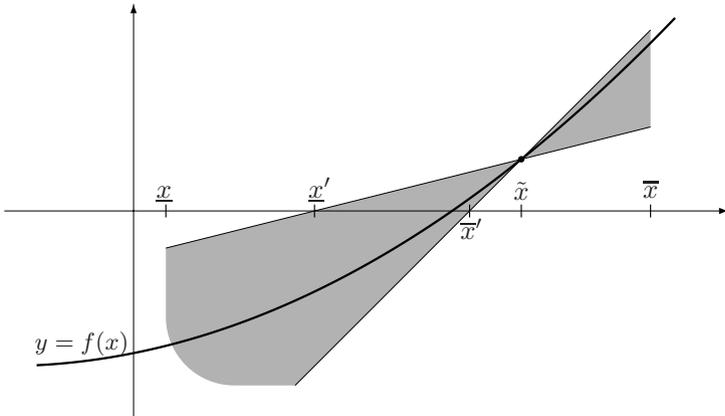


Рис. 4.9. Иллюстрация работы одномерного интервального метода Ньютона. Ситуация 1.

Допустим на время, что $0 \notin f'(x)$, так что $\mathcal{N}(x, \tilde{x})$ является вполне определённым конечным интервалом. Так как любой нуль функции $f(x)$ на x лежит также и в $\mathcal{N}(x, \tilde{x})$, то разумно взять в качестве следующего более точного приближения к решению пересечение

$$x \cap \mathcal{N}(x, \tilde{x}),$$

которое окажется, по крайней мере, не хуже x .

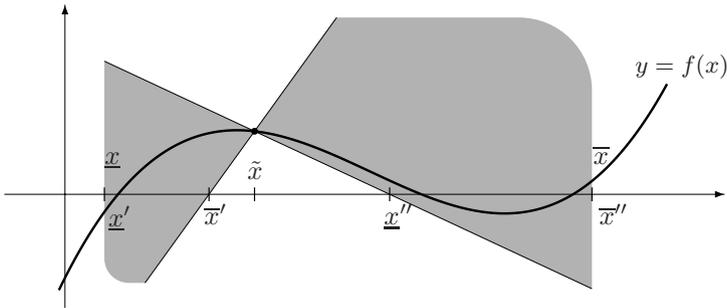


Рис. 4.10. Иллюстрация работы одномерного интервального метода Ньютона. Ситуация 2.

Далее, если $0 \in f'(x)$, мы можем придать смысл оператору Ньютона, воспользовавшись интервальной арифметикой Кахана. В действительности, эта модификация даже усилит интервальный метод Ньютона, так как мы получим возможность отделять решения друг от друга: в результате выполнения шага интервального метода Ньютона при $0 \in \text{int } f'(x)$ получаются, как правило, два непересекающихся интервала.

В арифметике Кахана дополнительно определено деление интервалов a и b с $0 \in b$, которое и приводит к бесконечным интервалам. Для удобства мы выпишем соответствующие результаты в развёрнутой

форме:

$$\begin{aligned}
 \underline{a}/\underline{b} &= \frac{[\underline{a}, \bar{a}]}{[\underline{b}, \bar{b}]} \\
 &= \begin{cases} \underline{a} \cdot [1/\bar{b}, 1/\underline{b}], & \text{если } 0 \notin \underline{b}, \\] - \infty, +\infty[, & \text{если } 0 \in \underline{a} \text{ и } 0 \in \underline{b}, \\ [\bar{a}/\underline{b}, +\infty[, & \text{если } \bar{a} < 0 \text{ и } \underline{b} < \bar{b} = 0, \\] - \infty, \bar{a}/\underline{b}] \cup [\bar{a}/\underline{b}, +\infty[, & \text{если } \bar{a} < 0 \text{ и } \underline{b} < 0 < \bar{b}, \\] - \infty, \bar{a}/\bar{b}], & \text{если } \bar{a} < 0 \text{ и } 0 = \underline{b} < \bar{b}, \\] - \infty, \underline{a}/\underline{b}], & \text{если } 0 < \underline{a} \text{ и } \underline{b} < \bar{b} = 0, \\] - \infty, \underline{a}/\underline{b}] \cup [\underline{a}/\bar{b}, +\infty[, & \text{если } 0 < \underline{a} \text{ и } \underline{b} < 0 < \bar{b}, \\ [\underline{a}/\bar{b}, +\infty[, & \text{если } 0 < \underline{a} \text{ и } 0 = \underline{b} < \bar{b}, \\ \emptyset, & \text{если } 0 \notin \underline{a} \text{ и } 0 = \underline{b}. \end{cases} \tag{4.24}
 \end{aligned}$$

Итак, перечислим основные свойства одномерного интервального метода Ньютона:

1. Всякий нуль функции f на исходном интервале \mathbf{x} корректно выделяется методом.
2. Если на исходном интервале \mathbf{x} нет нулей функции f , то этот факт будет установлен методом за конечное число итераций.
3. Если $0 \notin f'(\mathbf{x})$ для некоторого \mathbf{x} , то на следующем шаге метода будет исключена по крайней мере половина \mathbf{x}
4. Если $0 \notin f'(\mathbf{x})$, то асимптотический порядок сходимости метода к нулю функции f на интервале \mathbf{x} является квадратичным.

4.7в Многомерный интервальный метод Ньютона

Переходя к решению систем нелинейных уравнений, следует отметить, что многомерные версии интервального метода Ньютона гораздо

более многочисленны, чем одномерные, и отличаются очень большим разнообразием. В многомерном случае мы можем варьировать не только выбор точки \tilde{x} , вокруг которой осуществляется разложение, форму интервального расширения производных или наклонов функции, как это было в одномерном случае, но также и способ внешнего оценивания множества решений интервальной линейной системы, к которой приводится оценивание бруса решения. В оставшейся части этого параграфа мы рассмотрим простейшую форму многомерного интервального метода Ньютона, а его более специальными версиям, которые связываются с именами Кравчика и Хансена-Сенгупты, будут посвящены отдельные параграфы.

Определение 4.7.2 [46] *Для отображения $F : \mathbb{R}^n \supseteq D_0 \rightarrow \mathbb{R}^m$ матрица $A \in \mathbb{IR}^{m \times n}$ называется интервальной матрицей наклонов на $D \subseteq D_0$, если для любых $x, y \in D$ равенство*

$$F(y) - F(x) = A(y - x)$$

имеет место с некоторой вещественной $m \times n$ -матрицей $A \in \mathbf{A}$.

Предположим, что на брусе \mathbf{x} к решению предъявлена система нелинейных уравнений

$$F(x) = 0. \quad (4.25)$$

Если \mathbf{S} — интервальная матрица наклонов отображения F на \mathbf{x} , то для любых точек $x, \tilde{x} \in \mathbf{x}$ справедливо представление

$$F(x) \in F(\tilde{x}) + \mathbf{S}(x - \tilde{x}).$$

В частности, если x — решение системы уравнений (4.25), т. е. $F(x) = 0$, то

$$0 \in F(\tilde{x}) + \mathbf{S}(x - \tilde{x}). \quad (4.26)$$

Вспомним характеристику Бека для объединённого множества решений ИСЛАУ (Теорема 4.6.1): получается, что точка x удовлетворяет включению (4.26) тогда и только тогда, когда она принадлежит объединённому множеству решений интервальной линейной системы

$$\mathbf{S}(x - \tilde{x}) = -F(\tilde{x}). \quad (4.27)$$

Далее, если $Encl$ — процедура внешнего оценивания множества решений ИСЛАУ, то справедливо включение

$$x - \tilde{x} \in Encl(\mathbf{S}, -F(\tilde{x})),$$

так что

$$x \in \tilde{x} + \text{Encl}(\mathbf{S}, -F(\tilde{x})).$$

Определение 4.7.3 Пусть для внешнего оценивания множеств решений ИСЛАУ зафиксирована процедура Encl , а для отображения $F : \mathbb{R}^n \supseteq D \rightarrow \mathbb{R}^n$ известна интервальная матрица наклонов $\mathbf{S} \in \mathbb{IR}^{n \times n}$. Отображение

$$\mathcal{N} : \mathbb{ID} \times \mathbb{R}^n \rightarrow \mathbb{IR}^n,$$

задаваемое правилом

$$\mathcal{N}(\mathbf{x}, \tilde{x}) = \tilde{x} + \text{Encl}(\mathbf{S}, -F(\tilde{x})),$$

называется интервальным оператором Ньютона на \mathbb{ID} относительно точки \tilde{x} .

Как лучше выбирать центр разложения \tilde{x} ? Имеет смысл делать это так, чтобы величина $\|F(\tilde{x})\|$ была, по-возможности, меньшей. Чем меньше будет норма вектор-функции $F(\tilde{x})$, тем меньшим будет норма векторов, образующих множество решений интервальной линейной системы

$$\mathbf{S}(x - \tilde{x}) = -F(\tilde{x}),$$

которое мы должны пересекать с исходным брусом. Может быть, мы получим при этом более узкую внешнюю оценку множества решений исходной нелинейной системы и более точно определим статус исследуемого бруса. Численные эксперименты как будто подтверждают этот вывод.

Процедуру для уточнения центра разложения можно организовать как метод типа Ньютона, коль скоро нам известна интервальная матрица наклонов.

Наиболее неблагоприятной ситуацией при работе интервального метода Ньютона является, конечно, появление включения

$$\mathcal{N}(\mathbf{x}, \tilde{x}) \supseteq \mathbf{x}.$$

Тогда все последующие шаги зацикливаются на брусе \mathbf{x} и не дают никакой дополнительной информации об искомым решениях системы. Как поступать в этом случае? Ответ на этот вопрос рассматривается в следующем §4.8.

4.7г Метод Кравчика

Пусть на брус $\mathbf{x} \in \mathbb{IR}^n$ задана система n нелинейных уравнений с n неизвестными

$$F(x) = 0,$$

для которой требуется уточнить двусторонние границы решений. Возьмём какую-нибудь точку $\tilde{x} \in \mathbf{x}$ и организуем относительно неё разложение функции F :

$$F(x) \in F(\tilde{x}) + \mathbf{S}(x - \tilde{x}),$$

где $\mathbf{S} \in \mathbb{R}^{n \times n}$ — интервальная матрица наклонов отображения F на брус \mathbf{x} . Если x — это точка решения системы, то

$$0 \in F(\tilde{x}) + \mathbf{S}(x - \tilde{x}). \quad (4.26)$$

Но далее, в отличие от интервального метода Ньютона, мы не будем переходить к рассмотрению интервальной линейной системы (4.27), а домножим обе части этого включения слева на точечную $n \times n$ -матрицу, которую нам будет удобно обозначить как $(-A)$:

$$0 \in -AF(\tilde{x}) - A\mathbf{S}(x - \tilde{x}).$$

Добавление к обеим частям получившегося соотношения по $(x - \tilde{x})$ приводит к

$$x - \tilde{x} \in -AF(\tilde{x}) - A\mathbf{S}(x - \tilde{x}) + (x - \tilde{x}),$$

что равносильно

$$x \in \tilde{x} - AF(\tilde{x}) + (I - A\mathbf{S})(x - \tilde{x}),$$

так как для неинтервального общего множителя $(x - \tilde{x})$ можно воспользоваться дистрибутивным соотношением (1.16). Наконец, если решение x системы уравнений предполагается принадлежащим брусу \mathbf{x} , мы можем взять интервальное расширение по $x \in \mathbf{x}$ правой части полученного включения, придя к соотношению

$$x \in \tilde{x} - AF(\tilde{x}) + (I - A\mathbf{S})(\mathbf{x} - \tilde{x}),$$

Определение 4.7.4 Пусть определены некоторые правила, сопоставляющие всякому брусу $\mathbf{x} \in \mathbb{IR}^n$ точку $\tilde{x} \in \mathbf{x}$ и вещественную $n \times n$ -матрицу A и пусть также $\mathbf{S} \in \mathbb{IR}^{n \times n}$ — интервальная матрица наклонов отображения $F : \mathbb{R}^n \supseteq D \rightarrow \mathbb{R}^n$ на D . Отображение

$$\mathcal{K} : \mathbb{ID} \times \mathbb{R} \rightarrow \mathbb{IR}^n,$$

задаваемое выражением

$$\mathcal{K}(\mathbf{x}, \tilde{x}) := \tilde{x} - \Lambda F(\tilde{x}) + (I - \Lambda \mathbf{S})(\mathbf{x} - \tilde{x}),$$

называется оператором Кравчика на $\mathbb{I}D$ относительно точки \tilde{x} .

Теорема 4.7.2 Пусть $F : \mathbb{R}^n \supseteq D \rightarrow \mathbb{R}^n$ — непрерывное по Липшицу отображение, \mathbf{S} — его интервальная матрица наклонов и $\tilde{x} \in \mathbf{x} \subseteq \mathbb{I}D$. Тогда

- (i) каждое решение системы $F(x) = 0$ на брус \mathbf{x} лежит также в $\mathcal{K}(\mathbf{x}, \tilde{x})$;
- (ii) если $\mathbf{x} \cap \mathcal{K}(\mathbf{x}, \tilde{x}) = \emptyset$, то в \mathbf{x} нет решений системы $F(x) = 0$;
- (iii) если $\mathcal{K}(\mathbf{x}, \tilde{x}) \subseteq \mathbf{x}$, то в \mathbf{x} находится хотя бы одно решение системы $F(x) = 0$;
- (iv) если $\tilde{x} \in \text{int } \mathbf{x}$ и $\emptyset \neq \mathcal{K}(\mathbf{x}, \tilde{x}) \subseteq \text{int } \mathbf{x}$, то матрица \mathbf{S} сильно неособенна и в $\mathcal{K}(\mathbf{x}, \tilde{x})$ содержится в точности одно решение системы $F(x) = 0$.

Оператор Кравчика — это не что иное, как центрированная форма интервального расширения отображения $\Phi(x) = x - \Lambda F(x)$, возникающего в правой части системы уравнений после её приведения к рекуррентному виду

$$x = \Phi(x).$$

4.8 Глобальное решение уравнений и систем уравнений

Если ширина бруса \mathbf{X} велика, то на нём описанные в предшествующем параграфе методики уточнения решения могут оказаться малоуспешными в том смысле, что мы получим включение (4.23), из которого нельзя вывести никакого определённого заключения ни о существовании решения на брус \mathbf{X} , ни о его отсутствии. Кроме того, сам этот брус, как область потенциально содержащая решение, несколько не будет уточнён (уменьшен).

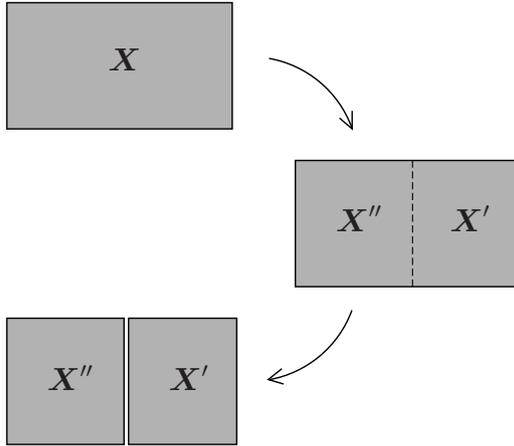


Рис. 4.11. Принудительное дробление бруса.

Тогда практикуют принудительное дробление \mathbf{X} на более мелкие подбрусы. Наиболее популярна при этом *бисекция* — разбиение бруса \mathbf{X} на две (равные или неравные) части вдоль какой-нибудь грани, например, на половинки

$$\begin{aligned}\mathbf{X}' &= (\mathbf{X}_1, \dots, [\underline{\mathbf{X}}_\iota, \text{mid } \mathbf{X}_\iota], \dots, \mathbf{X}_n), \\ \mathbf{X}'' &= (\mathbf{X}_1, \dots, [\text{mid } \mathbf{X}_\iota, \overline{\mathbf{X}}_\iota], \dots, \mathbf{X}_n)\end{aligned}$$

для некоторого номера $\iota \in \{1, 2, \dots, n\}$. При этом подбрусы \mathbf{X}' и \mathbf{X}'' называются *потомками* бруса \mathbf{X} . Далее эти потомки можно разбить ещё раз, и ещё \dots — столько, сколько необходимо для достижения желаемой малости их размеров, при которой мы сможем успешно выполнять на этих брусах рассмотренные выше тесты существования решений.

Если мы не хотим упустить при этом ни одного решения системы, то должны хранить все возникающие в процессе такого дробления подбрусы, относительно которых тестом существования не доказано строго, что они не содержат решений. Организуем поэтому *рабочий список* \mathcal{L} из всех потомков начального бруса \mathbf{X} , подозрительных на содержание решений. Хотя мы называем эту структуру данных «списком», в смысле программной реализации это может быть не обязательно список, а любое хранилище брусов, организованное, к примеру, как *стек*

(магазин) или *куча* и т. п. (см. [4]) В целом же алгоритм глобального доказательного решения системы уравнений организуем в виде повторяющейся последовательности следующих действий:

- извлечение некоторого бруса из списка \mathcal{L} ,
- дробление этого бруса на потомки,
- проверка существования решений в каждом из подбрус-потомков, по результатам которой мы
 - либо выдаём этот подбрус в качестве ответа к решаемой задаче,
 - либо заносим его в рабочий список \mathcal{L} для последующей обработки,
 - либо исключаем из дальнейшего рассмотрения, как не содержащий решений рассматриваемой системы.

Кроме того, чтобы обеспечить ограниченность времени работы алгоритма, на практике имеет смысл задаться некоторым порогом мелкости (малости размеров) брусков δ , при достижении которого дальше дробить брус уже не имеет смысла. В Табл. 4.2 приведён псевдокод получающегося алгоритма, который называется *методом ветвлений и отсечений*: ветвления соответствуют разбиениям исходного бруса на подбрусы (фактически, разбиениям исходной задачи на подзадачи), а отсечения — это отбрасывание бесперспективных подбрусков исходной области поиска.²

Отметим, что неизбежные ограничения на вычислительные ресурсы ЭВМ могут воспрепятствовать решению этим алгоритмом задачи (4.19) «до конца», поскольку могут возникнуть ситуации, когда

- 1) размеры обрабатываемого бруса уже меньше δ , но нам ещё не удаётся ни доказать существование в нём решений, ни показать их отсутствие;
- 2) размеры обрабатываемого бруса всё ещё больше δ , но вычислительные ресурсы уже не позволяют производить его обработку дальше: исчерпались выделенное время, память и т. п.

²Стандартный английский термин для обозначения подобного типа алгоритмов — «branch-and-prune». С ними тесно связаны *методы ветвей и границы*, широко применяемые в вычислительной оптимизации.

Таблица 4.2. Интервальный метод ветвлений и отсечений для глобального доказательного решения уравнений

<p>Вход</p> <p>Система уравнений $F(x) = 0$. Брус $\mathbf{X} \in \mathbb{IR}^n$. Интервальное расширение $\mathbf{F} : \mathbb{IX} \rightarrow \mathbb{IR}^n$ функции F. Заданная точность $\delta > 0$ локализации решений системы.</p>
<p>Выход</p> <p>Список НавернякаРешения из брусов размера менее δ, которые гарантированно содержат решения системы уравнений в \mathbf{X}. Список ВозможноРешения из брусов размера менее δ, которые могут содержать решения системы уравнений в \mathbf{X}. Список Недообработанные из брусов размера более δ, которые могут содержать решения системы уравнений в \mathbf{X}.</p>
<p>Алгоритм</p> <p>инициализируем рабочий список \mathcal{L} исходным брусом \mathbf{X} ; DO WHILE (($\mathcal{L} \neq \emptyset$) и (не исчерпаны ресурсы ЭВМ)) извлекаем из рабочего списка \mathcal{L} брус \mathbf{Y} ; применяем к \mathbf{Y} тест существования решения, его результат обозначаем также через \mathbf{Y} ; IF (в \mathbf{Y} доказано отсутствие решений) THEN удаляем брус \mathbf{Y} из рассмотрения ELSE IF ((размер бруса \mathbf{Y}) $< \delta$) THEN вносим \mathbf{Y} в соответствующий из списков НавернякаРешения или ВозможноРешения ELSE рассекаем \mathbf{Y} на потомки \mathbf{Y}' и \mathbf{Y}'' и вносим их в рабочий список \mathcal{L} END IF END IF END DO все брусы из \mathcal{L} перемещаем в список Недообработанные;</p>

В реальных вычислениях остановка алгоритма Табл. 4.2 может происходить поэтому не только при достижении пустого рабочего списка \mathcal{L} (когда исчерпана вся область поиска решений), но и, к примеру, при достижении определённого числа шагов или времени счёта и т. п. Тогда все брусы, оставшиеся в рабочем списке \mathcal{L} , оказываются не до конца обработанными, и мы условимся так и называть их — «недообработанные». Итак, в общем случае результатом работы нашего алгоритма должны быть три списка брусков:

список **НавернякаРешения**, состоящий из брусков шириной меньше δ , которые гарантированно содержат решения,

список **ВозможноРешения**, состоящий из брусков шириной меньше δ , подозрительных на содержание решения, и

список **Недообработанные**, состоящий брусков, которые алгоритму не удалось обработать «до конца» и которые имеют ширину не меньше δ .

При этом все решения рассматриваемой системы уравнений, не принадлежащие брусам из списка **НавернякаРешения**, содержатся в брусах из списков **ВозможноРешения** и **Недообработанные**.

Практика эксплуатации интервальных методов для доказательно-глобального решения уравнений и систем уравнений выявила ряд проблем и трудностей. Во многих случаях (особенно при наличии так называемых кратных корней) задачу не удаётся решить до конца и предъявить все гарантированные решения уравнения. Список брусков-ответов с неопределённым статусом (**ВозможноРешения** в псевдокоде Табл. 4.2) часто никак не собираются исчезать ни при увеличении точности вычислений, ни при выделении дополнительного времени счёта и т.п. Нередко он разрастается до огромных размеров, хотя большинство образующих его брусков возможных решений являются «фантомами» немногих реальных решений. Но эти феномены могут быть успешно объяснены на основе теории, изложенной в §4.3.

Решения уравнений и систем уравнений — это особые точки соответствующих векторных полей, которые, как мы могли видеть, отличаются большим разнообразием. Насколько используемые нами при доказательном решении систем уравнений инструменты приспособлены для выявления особых точек различных типов? Нетрудно понять, что интервальный метод Ньютона, методы Кравчика и Хансена-Сенгупты,

тесты существования Мура и Куи, основывающиеся на теоремах Лерэ-Шаудера и Брауэра и наиболее часто используемые при практических доказательных вычислениях решений уравнений, охватывают только случаи индекса ± 1 особой точки F . Если же решение системы является критической точкой соответствующего отображения с индексом, не равным ± 1 , то доказать его существование с помощью вышеупомянутых результатов принципиально не получится. Это объясняет, почему многие существующие практические интервальные алгоритмы для доказательного глобального решения систем уравнений не могут достичь «полного успеха» в общем случае.

Помимо вышеназванной причины необходимо отметить, что список **ВозможноРешения** может соответствовать неустойчивым решениям системы уравнений, имеющим нулевой индекс. Эти решения разрушаются при сколь угодно малых возмущениях уравнений и потому не могут быть идентифицированы никаким приближенным вычислительным алгоритмом с конечной точностью представления данных. К примеру, таковым является кратный корень квадратного уравнения (4.5)–(4.6), и хорошо известно, что он плохо находится численно как традиционными, так и интервальными подходами.

Алгоритмы ветвлений и отсечений, дополненные различными усовершенствованиями и приёмами, ускоряющими сходимость, получили большое развитие в интервальном анализе в последние десятилетия (см., например, книги [35, 41, 45, 46]), а реализованные на их основе программные комплексы существенно продвинули практику численного решения уравнений и систем уравнений.

Литература к главе 4

Основная

- [1] АЛЕФЕЛЬД Г., ХЕРЦБЕРГЕР Ю. *Введение в интервальные вычисления*. – Москва: Мир, 1987.
- [2] АКРИТАС А. *Основы компьютерной алгебры с приложениями*. – Москва: Мир, 1994.
- [3] БАРАХНИН В.Б., ШАПЕЕВ В.П. *Введение в численный анализ*. – Санкт-Петербург–Москва–Краснодар: Лань, 2005.
- [4] БАУЭР Ф.Л., ГООЗ Г. *Информатика. В 2-х ч.* – Москва: Мир, 1990.
- [5] БАХВАЛОВ Н.С., ЖИДКОВ Н.П., КОБЕЛЬКОВ Г.М. *Численные методы*. – Москва: Бином, 2003, а также другие издания этой книги.

- [6] Бахвалов Н.С., Корнев А.А., Чижонков Е.В. *Численные методы. Решение задач и упражнения.* – Москва: Дрофа, 2008.
- [7] Березин И.С., Жидков Н.П. *Методы вычислений. Т. 1–2.* – Москва: Наука, 1966.
- [8] Берже М. *Геометрия. Т. 1, 2.* – Москва: Наука, 1984.
- [9] Бержвицкий В.М. *Численные методы. Части 1–2.* – Москва: «Оникс 21 век», 2005.
- [10] Волков Е.А. *Численные методы.* – Москва: Наука, 1987.
- [11] Годунов С.К. *Современные аспекты линейной алгебры.* – Новосибирск: Научная книга, 1997.
- [12] Годунов С.К., Антонов А.Г., Кириллюк О.П., Костин В.И. *Гарантированная точность решения систем линейных уравнений в евклидовых пространствах.* – Новосибирск: Наука, 1992.
- [13] Гэри М., Джонсон Д. *Вычислительные машины и труднорешаемые задачи.* – Москва: Мир, 1982.
- [14] Демидович Б.П., Марон А.А. *Основы вычислительной математики.* – Москва: Наука, 1970.
- [15] Дэннис Дж., мл., Шнабель Р. *Численные методы безусловной оптимизации и решения нелинейных уравнений.* – Москва: Мир, 1988.
- [16] Калиткин Н.Н. *Численные методы.* – Москва: Наука, 1978.
- [17] Канторович Л.В., Акилов Г.П. *Функциональный анализ.* – Москва: Наука, 1984.
- [18] Коллатц Л. *Функциональный анализ и вычислительная математика.* – Москва: Мир, 1969.
- [19] Крылов А.Н. *Лекции о приближённых вычислениях.* – Москва: ГИТТЛ, 1954, а также более ранние издания.
- [20] Крылов В.И., Бобков В.В., Монастырный П.И. *Вычислительные методы. Т. 1–2.* – Москва: Наука, 1976.
- [21] Кунц К.С. *Численный анализ.* – Киев: Техника, 1964.
- [22] Мацокин А.М. *Численный анализ. Вычислительные методы линейной алгебры. Конспекты лекций для преподавания в III семестре ММФ НГУ.* – Новосибирск: НГУ, 2009–2010.
- [23] Мацокин А.М., Сорокин С.Б. *Численные методы. Часть 1. Численный анализ.* – Новосибирск: НГУ, 2006.
- [24] Меньшиков Г.Г. *Локализуемые вычисления. Конспект лекций.* – Санкт-Петербург: СПбГУ, Факультет прикладной математики–процессов управления, 2003.
- [25] Миньков С.Л., Миньков Л.Л. *Основы численных методов.* – Томск: Издательство научно-технической литературы, 2005.
- [26] Мысовских И.П. *Лекции по методам вычислений.* – Санкт-Петербург: Издательство Санкт-Петербургского университета, 1998.

- [27] ОРТЕГА Дж., РЕЙНБОЛДТ В. *Итерационные методы решения нелинейных систем уравнений со многими неизвестными*. – Москва: Мир, 1975.
- [28] ОСТРОВСКИЙ А.М. *Решение уравнений и систем уравнений*. – Москва: Издательство иностранной литературы, 1963.
- [29] САМАРСКИЙ А.А., ГУЛИН А.В. *Численные методы*. – Москва: Наука, 1989.
- [30] СЕМЁНОВ А.Л., ВАЖЕВ И.В., КАШЕВАРОВА Т.П. и др. Интервальные методы распространения ограничений и их приложения // *Системная информатика*. – Новосибирск: Издательство СО РАН, 2004. – Вып. 9. – С. 245–358.
- [31] ТРАВУ ДЖ. *Итерационные методы решения уравнений*. – Москва: Мир, 1985.
- [32] ТЫРТЫШНИКОВ Е.Е. *Методы численного анализа*. – Москва: Академия, 2007.
- [33] УСПЕНСКИЙ В.А., СЕМЁНОВ А.Л. *Теория алгоритмов: основные открытия и приложения*. – Москва: Наука, 1987.
- [34] ФИХТЕНГОЛЬЦ Г.М. *Курс дифференциального и интегрального исчисления. Т. 1*. – Москва: Наука, 1966.
- [35] ХАНСЕН Э., УОЛСТЕР ДЖ.У. *Глобальная оптимизация с помощью методов интервального анализа*. – Москва-Ижевск: Издательство «РХД», 2012.
- [36] ХОЛОДНИОК М., КЛИЧ А., КУБИЧЕК М., МАРЕК М. *Методы анализа нелинейных динамических моделей*. – Москва: Мир, 1991.
- [37] ШАРЫЙ С.П. *Конечномерный интервальный анализ*. – Электронная книга, 2010 (см. <http://www.nsc.ru/interval/Library/InteBooks>)
- [38] ШИЛОВ Г.Е. *Математический анализ. Функции одного переменного. Ч. 1–2*. – Москва: Наука, 1969.
- [39] АВЕРТН О. *Precise numerical methods using C++*. – San Diego: Academic Press, 1998.
- [40] AKYILDIZ Y., AL-SUWAIYEL M.I. No pathologies for interval Newton's method // *Interval Computations*. – 1993. – No. 1. – P. 60–72.
- [41] KEARFOTT R.B. *Rigorous global search: Continuous problems*. – Dordrecht: Kluwer, 1996.
- [42] KELLEY C.T. *Iterative methods for linear and nonlinear equations*. – Philadelphia: SIAM, 1995.
- [43] KREINOVICH V., LAKEYEV A.V., ROHN J., KAHL P. *Computational complexity and feasibility of data processing and interval computations*. – Dordrecht: Kluwer, 1997.
- [44] MIRANDA C. Un' osservazione su un teorema di Brouwer // *Bollet. Unione Mat. Ital. Serie II*. – 1940. – Т. 3. – С. 5–7.
- [45] MOORE R.E., KEARFOTT R.B., CLOUD M. *Introduction to interval analysis*. – Philadelphia: SIAM, 2009.
- [46] NEUMAIER A. *Interval methods for systems of equations*. – Cambridge: Cambridge University Press, 1990.

- [47] TREFETHEN L.N. Pseudospectra of linear operators // SIAM Review. 1997. – Vol. 39, No. 3. – P. 383–406.
- [48] TREFETHEN L.N., BAU D. III *Numerical linear algebra*. – Philadelphia: SIAM, 1997.

Дополнительная

- [49] АБАФФИ Й., СПЕДИКАТО Э. *Математические методы для линейных и нелинейных уравнений. Проекционные ABS-алгоритмы*. – Москва: Мир, 1996.
- [50] АРНОЛЬД В.И. *Обыкновенные дифференциальные уравнения*. – Москва: Наука, 1984.
- [51] БАБЕНКО К.И. *Основы численного анализа*. – Москва: Наука, 1986.
- [52] ГАНШИН Г.С. *Методы оптимизации и решение уравнений*. – Москва: Наука, 1987.
- [53] ЗАГУСКИН В.Л. *Справочник по численным методам решения алгебраических и трансцендентных уравнений*. – Москва: Физматгиз, 1960.
- [54] КРАСНОСЕЛЬСКИЙ М.А., ЗАВРЕЙКО П.П. *Геометрические методы нелинейного анализа*. – Москва: Наука, 1975.
- [55] КРАСНОСЕЛЬСКИЙ М.А., ПЕРОВ А.И., ПОВОЛОЦКИЙ А.И., ЗАВРЕЙКО П.П. *Векторные поля на плоскости*. – Москва: Физматлит, 1963.
- [56] НИРЕНБЕРГ Л. *Лекции по нелинейному функциональному анализу*. – Москва: Мир, 1977.
- [57] ОПОЙЦЕВ В.И. *Нелинейная системостатика*. – Москва: Наука, 1986.
- [58] The NIST reference on constants, units, and uncertainty. – <http://physics.nist.gov/cuu/Constants>
- [59] Pseudospectra gateway. – <http://web.comlab.ox.ac.uk/projects/pseudospectra/>
- [60] Scilab — The Free Platform for Numerical Computation. <http://www.scilab.org>

Обозначения

\Rightarrow	логическая импликация
\Leftrightarrow	логическая равносильность
$\&$	логическая конъюнкция, связка «и»
\rightarrow	отображение множеств или предельный переход
\mapsto	правило сопоставления элементов при отображении
\leftarrow	оператор присваивания в алгоритмах
\circ	знак композиции отображений
\emptyset	пустое множество
$x \in X$	элемент x принадлежит множеству X
$x \notin X$	элемент x не принадлежит множеству X
$X \cup Y$	объединение множеств X и Y
$X \cap Y$	пересечение множеств X и Y
$X \setminus Y$	разность множеств X и Y
$X \subseteq Y$	множество X есть подмножество множества Y
$X \times Y$	прямое декартово произведение множеств X и Y
\mathbb{N}	множество натуральных чисел
\mathbb{R}	множество вещественных (действительных) чисел
\mathbb{R}_+	множество неотрицательных вещественных чисел
\mathbb{C}	множество комплексных чисел
$\mathbb{I}\mathbb{R}$	множество интервалов вещественной оси \mathbb{R}
\mathbb{R}^n	множество вещественных n -мерных векторов
\mathbb{C}^n	множество комплексных n -векторов

\mathbb{R}^n	множество n -мерных интервальных векторов
$\mathbb{R}^{m \times n}$	множество вещественных $m \times n$ -матриц
$\mathbb{C}^{m \times n}$	множество комплексных $m \times n$ -матриц
$\mathbb{R}^{m \times n}$	множество интервальных $m \times n$ -матриц
\approx	приблизительно равно
$\approx \leq$	приблизительно меньше или равно
i	мнимая единица
\bar{z}	комплексно сопряжённое к числу $z \in \mathbb{C}$
$\operatorname{sgn} a$	знак числа $a \in \mathbb{R}$
$[a, b]$	интервал с нижним концом a и верхним b
$]a, b[$	открытый интервал с концами a и b
$\underline{a}, \inf \mathbf{a}$	левый конец интервала \mathbf{a}
$\bar{a}, \sup \mathbf{a}$	правый конец интервала \mathbf{a}
$\operatorname{mid} \mathbf{a}$	середина интервала \mathbf{a}
$\operatorname{wid} \mathbf{a}$	ширина интервала \mathbf{a}
dist	метрика (расстояние)
Dist	мультиметрика (векторнозначное расстояние)
$f(x) _a^b$	разность значений функции f между $x = a$ и $x = b$
$\operatorname{dom} f$	область определения функции f
$\operatorname{ran}(f, X)$	область значений функции f на X
$f^\sphericalangle(\cdot)$	разделённая разность от функции f
\min, \max	операции взятия минимума и максимума
$\operatorname{int} X$	топологическая внутренность множества X
$\operatorname{cl} X$	топологическое замыкание множества X
∂X	граница множества X
δ_{ij}	символ Кронекера, 1 при $i = j$ и 0 иначе
I	единичная матрица соответствующих размеров
$\ \cdot\ $	векторная или матричная норма
$\langle \cdot, \cdot \rangle$	скалярное произведение векторов
A^\top	матрица, транспонированная к матрице A

A^*	матрица, эрмитово сопряжённая к матрице A
A^{-1}	матрица, обратная к матрице A
$\rho(A)$	спектральный радиус матрицы A
$\lambda(A), \lambda_i(A)$	собственные числа матрицы A
$\sigma(A), \sigma_i(A)$	сингулярные числа матрицы A
$\text{cond}(A)$	число обусловленности матрицы A
$\text{rank } A$	ранг матрицы A
$\det A$	определитель матрицы A
$\mathcal{K}_i(A, r)$	подпространство Крылова матрицы A
$\text{diag} \{z_1, \dots, z_n\}$	диагональная $n \times n$ -матрица с элементами z_1, \dots, z_n по главной диагонали
$\text{lin} \{v_1, \dots, v_n\}$	линейная оболочка векторов v_1, \dots, v_n
$C^p[a, b]$	класс функций, непрерывно дифференцируемых вплоть до p -го порядка на интервале $[a, b]$
$\mathcal{L}^2[a, b]$	класс функций, интегрируемых с квадратом на интервале $[a, b]$
\sum	символ суммы нескольких слагаемых
\prod	символ произведения нескольких сомножителей

Интервалы и другие интервальные величины (векторы, матрицы и др.) всюду в тексте обозначаются жирным математическим шрифтом, например, $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots, \mathbf{x}, \mathbf{y}, \mathbf{z}$, тогда как неинтервальные (точечные) величины никак специально не выделяются. Арифметические операции с интервальными величинами — это операции классической интервальной арифметики \mathbb{IR} (см. §1.4).

Если не оговорено противное, под векторами (точечными или интервальными) всюду понимаются вектор-столбцы.

Конец доказательства теоремы или предложения и конец примера выделяются в тексте стандартным знаком «■».

Значительная часть описываемых в книге алгоритмов снабжается псевдокодами на неформальном алгоритмическом языке, в котором операторные скобки

DO FOR ... END DO означают оператор цикла со счётчиком, который задаётся после FOR,

DO WHILE ... END DO означают оператор цикла с предусловием, стоящим после WHILE,

IF ... THEN ... END IF или IF ... THEN ... ELSE ... END IF означают условные операторы с условием, стоящим после IF.

В циклах DO FOR ключевое слово TO означает увеличение счётчика итераций от начального значения до конечного (положительный шаг), а ключевое слово DOWNTO — уменьшение счётчика итераций (отрицательный шаг).

Краткий биографический словарь

Абель, Нильс Хенрик (Niels Henrik Abel, 1802–1829)

— норвежский математик.

Адамар, Жак Саломон (Jacques Salomon Hadamard, 1865–1963)

— французский математик.

Андронов, Александр Александрович (1901–1952)

— советский физик и механик.

Бабенко, Константин Иванович (1919–1987)

— советский математик и механик.

Банах, Стефан (Stefan Banach, 1892–1945)

— польский математик.

Бауэр, Фридрих Людвиг (Friedrich Ludwig Bauer, род. 1924)

— немецкий математик.

Бельтрами, Эудженио (Eugenio Beltrami, 1835–1900)

— итальянский математик.

Бернштейн, Сергей Натанович (1880–1968)

— российский и советский математик.

Больцано, Бернارد (Bernard Bolzano, 1781–1848)

— чешский теолог, философ и математик.

Брадис, Владимир Модестович (1890–1975)

— русский и советский математик и педагог.

Брауэр, Лейтзен Эгберт Ян (Luitzen Egbertus Jan Brouwer, 1881–1966)

— голландский математик.

- Бюффон, Жорж-Луи Леклерк де (Georges-Louis Leclerc de Buffon, 1707–1788) — французский естествоиспытатель.
- Валлис, Джон (John Wallis, 1616–1703)
— английский математик.
- Вандермонд, Александр Теофиль (Alexandre Theophill Vandermonde, 1735–1796) — французский музыкант и математик.
- Вейерштрасс, Карл Теодор (Karl Theodor Weierstrass, 1815–1897)
— немецкий математик.
- Вейль, Герман (Hermann Weyl, 1885–1955)
— немецкий и американский математик.
- Виет, Франсуа (François Viète, 1540–1603)
— французский математик.
- Виландт, Хельмут (Helmut Wielandt, 1910–2001)
— немецкий математик.
- Гаусс, Карл Фридрих (Carl Friedrich Gauss, 1777–1855)
— немецкий математик, внёсший также фундаментальный вклад в численные методы, астрономию и геодезию.
- Гельфанд, Израиль Моисеевич (1913–2009)
— советский математик. С 1989 года жил и работал в США.
- Герон Александрийский (др.-греч. *Ἡρώων ο Αλεξανδρεὺς*, около 1 в. н.э.)
— греческий математик и механик.
- Гершгорин, Семён Аронович (1901–1933)
— советский математик, живший и работавший в Ленинграде.
- Гёльдер, Людвиг Отто (Ludwig Otto Hölder, 1859–1937)
— немецкий математик.
- Гивенс, Джеймс Уоллес (James Wallace Givens, 1910–1993)
— американский математик.
- Гильберт, Давид (David Hilbert, 1862–1943)
— немецкий математик.
- Грам, Йорген Педерсен (Jorgen Pedersen Gram, 1850–1916)
— датский математик.
- Евклид, или Эвклид (др.-греч. *Ευκλείδης*, около 300 г. до н. э.)
— древнегреческий математик.

Жордан, Мари Энмон Камилл (Marie Ennemond Camille Jordan, 1838–1922) — французский математик.

Зейдель, Филипп Людвиг (Philipp Ludwig Seidel, 1821–1896)
— немецкий астроном и математик.

Йордан, Вильгельм (Wilhelm Jordan, 1842–1899)
— немецкий геодезист.³

Канторович, Леонид Витальевич (1912–1986)
— советский математик и экономист, известный пионерским вкладом в линейное программирование.

Кнут, Дональд Эрвин (Donald Ervin Knuth, род. 1938)
— американский математик и специалист по информатике и программированию.

Колмогоров, Андрей Николаевич (1903–1987)
— советский математик, внёсший большой вклад во многие разделы современной математики, от топологии до теории вероятностей.

Котес, Роджер (Roger Cotes, 1682–1716)
— английский математик.

Коши, Огюстен Луи (Augustin Louis Cauchy, 1789–1857)
— французский математик и механик.

Кравчик, Рудольф (Krawczyk, Rudolf)
— немецкий математик.

Крамер, Габриэль (Gabriel Cramer, 1704–1752)
— швейцарский математик.

Красносельский, Марк Александрович (1920–1997)
— советский и российский математик.

Крейн, Селим Григорьевич (1917–1999)
— советский и российский математик.

Кронекер, Леопольд (Leopold Kronecker, 1823–1891)
— немецкий математик.

Крылов, Алексей Николаевич (1863–1945)
— русский и советский математик, механик и кораблестроитель.

Кублановская, Вера Николаевна (род. 1920)
— советский и российский математик.

³Не следует путать его с Паскуалем Йорданом (Pascual Jordan, 1902–1980), немецким физиком и математиком.

- Кузьмин, Родион Осиевич (1891–1949)
— русский и советский математик.
- Курант, Рихард (Richard Courant, 1888–1972)
— немецкий и американский математик.
- Лагранж, Жозеф Луи (Joseph Louis Lagrange, 1736–1813)
— французский математик и механик.
- Ландау, Эдмунд (Edmund Landau, 1877–1938)
— немецкий математик.
- Ланцош, Корнелий (Cornelius Lanczos, 1893–1974)
— американский физик и математик венгерского происхождения.
- Пьер-Симон, Лаплас (Pierre-Simon Laplace, 1749–1827)
— французский математик, механик, физик и астроном.
- Лебег, Анри Леон (Henri Léon Lebesgue, 1875–1941)
— французский математик.
- Лежандр, Адриен Мари (Adrien Marie Legendre, 1752–1833)
— французский математик и механик.
- Лейбниц, Готфрид Вильгельм (Gottfried Wilhelm Leibnitz, 1646–1716)
— немецкий философ, математик и физик, один из создателей дифференциального и интегрального исчисления.
- Липшиц, Рудольф (Rudolf Lipschitz, 1832–1903)
— немецкий математик.
- Лобачевский, Николай Иванович (1792–1856)
— русский математик, создатель неевклидовой геометрии.
- Локуциевский, Олег Вячеславович (1922–1990)
— советский математик.
- Ляпунов, Александр Михайлович (1857–1918)
— русский математик и механик, основоположник математической теории устойчивости.
- Марков, Андрей Андреевич (1856–1922)
— русский математик.
- Марцинкевич, Юзеф (Józef Marcinkiewicz, 1910–1941)
— польский математик.
- Микеладзе, Шалва Ефимович (1895–1976)
— советский математик.

- Минковский, Герман (Hermann Minkowski, 1864–1909)
— немецкий математик.
- Миранда, Карло (Carlo Miranda, 1912–1982)
— итальянский математик.
- Нейман, Карл Готфрид (Karl Gottfried Neumann, 1832–1925)
— немецкий математик.
- фон Нейман, Джон (John von Neumann, 1903–1957)
— американский математик венгерского происхождения.⁴
- Ньютон, Исаак (Isaac Newton, 1643–1727)
— английский физик и математик, заложивший основы дифференциального и интегрального исчисления и механики.
- Островский, Александр (Alexander M. Ostrowski, 1893–1986)
немецкий и швейцарский математик русского происхождения.
- Перрон, Оскар (Oskar Perron, 1880–1975)
— немецкий математик.
- Пикар, Шарль Эмиль (Picard, Charles Émile, 1856–1941)
— французский математик.
- Пирсон, Карл (Чарльз) (Karl (Charles) Pearson, 1857–1936)
— английский математик, биолог и философ.
- Пойа (Полия), Дьёрдь (иногда Джордж) (György Pólya, 1887–1985)
— венгерский и американский математик.
- Риман, Бернхард (Georg-Friedrich-Bernhard Riemann, 1826–1866)
— немецкий математик, механик и физик.
- Ричардсон, Льюис Фрай (Lewis Fry Richardson, 1881–1953)
— английский математик, физик и метеоролог.
- Родриг, Бенжамен Оленд (Benjamin Olinde Rodrigues, 1795–1851)
— французский математик и банкир.
- Рунге, Карл Давид (Karl David Runge, 1856–1927)
— немецкий физик и математик.
- Руффини, Паоло (Paolo Ruffini, 1765–1822)
— итальянский математик.

⁴Его именем, в частности, назван спектральный признак устойчивости разностных схем.

- Рэлей, Джон Уильям (John William Reyleigh, 1842–1919)
— английский физик.
- Самарский, Александр Андреевич (1919–2008)
— советский и российский математик.
- Симпсон, Томас (Thomas Simpson, 1710–1761)
— английский математик.
- Сонин Николай Яковлевич (1849–1915)
— русский математик.
- Стеклов, Владимир Андреевич (1863–1926)
— русский математик и механик.
- Стирлинг, Джеймс (James Stirling, 1692–1770)
— шотландский математик.
- Таусски, Ольга (Olga Tausski, 1906–1995)
— американский математик.
- Тейлор, Брук (Brook Taylor, 1685–1731)
— английский математик.
- Тихонов, Андрей Николаевич (1906–1993)
— советский математик.
- Улам, Станислав (Stanislaw Marcin Ulam, 1909–1984)
американский математик польского происхождения.
- Фабер, Георг (Georg Faber, 1877–1966)
— немецкий математик.
- Фаддеев, Дмитрий Константинович (1907–1989)
— советский математик.
- Фаддеева, Вера Николаевна (1906–1983)
— советский математик.
- Файк, (С.Т. Fike, –)
— американский математик.
- Фарадей, Майкл (Michael Faraday, 1791–1867)
— английский физик и химик.
- Федоренко, Радий Петрович (1930–2009)
— советский математик.
- Ферма, Пьер (Pierre Fermat, 1601–1665)
— французский математик.

- Фишер, Эрнст Сигизмунд (Ernst Sigismund Fischer, 1875–1954)
— немецкий математик.⁵
- Фробениус, Фердинанд Георг (Ferdinand Georg Frobenius, 1849–1917)
— немецкий математик.
- Фрэнсис, Джон (John G.F. Francis, род. 1934)
— английский математик и программист.
- Хаусдорф, Феликс (Felix Hausdorff, 1868–1942)
— немецкий математик.
- Хаусхолдер, Элстон (Alston Scott Householder, 1904–1993)
— американский математик.
- Хессенберг, Карл Адольф (Karl Adolf Hessenberg, 1904–1959)
— немецкий математик и инженер.
- Хестенс, Магнус (Magnus R. Hestenes, 1906–1991)
— американский математик.
- Холлесский, Андре-Луи (André-Louis Cholesky, 1875–1918)
— французский геодезист и математик.⁶
- Хопф, Хайнц (Heinz Hopf, 1896–1971)
— немецкий и швейцарский математик.
- Хоффман, Алан Джером (Alan Jerome Hoffman, род. 1924)
— американский математик.⁷
- Чебышёв, Пафнутий Львович (1821–1894)
— русский математик и механик, внёсший основополагающий вклад,
в частности, в теорию приближений и теорию вероятностей.
- Шёнберг, Исаак Якоб (Isaac Jacob Schönberg, 1903–1990)
румынский и американский математик.
- Шмидт, Эрхард (Erhard Schmidt, 1876–1959)
— немецкий математик.
- Шрёдер, Иоганн (Johann Schröder, 1925–2007)
— немецкий математик.

⁵К этому же времени относится жизнь и деятельность Рональда Э. Фишера (1890–1962), английского статистика, с именем которого связаны важные результаты математической статистики.

⁶В русской научной литературе его фамилия нередко транслитерируется как «Холецкий» или даже «Халецкий».

⁷Иногда его фамилию транслитерируют как «Гоффман».

Штифель, Эдуард (Eduard L. Stiefel, 1909–1978)
швейцарский математик.

Шур, Исай (Issai Schur, 1875–1941)
— немецкий и израильский математик.

Эйлер, Леонард (Leonhard Euler, 1707–1783)
— российский математик швейцарского происхождения, внёсший
фундаментальный вклад практически во все разделы математики.

Эрмит, Шарль (Charles Hermite, 1822–1901)
— французский математик.

Якоби, Карл Густав (Carl Gustav Jacobi, 1804–1851)
— немецкий математик.

Яненко, Николай Николаевич (1921–1984)
— советский математик и механик.

Предметный указатель

- ε -решения, 441
- p -ранговое приближение матрицы, 257
- абсолютная погрешность, 11
- алгебраическая степень точности, 145
- алгоритмическое дифференцирование, 99, 118
- аналитическая функция, 82
- арифметика дифференциальная, 118
- автоматическое дифференцирование, 99, 118
- биортогональность, 210
- целевая функция, 359
- чебышёвская метрика, 43
- чебышёвская норма, 226
- чебышёвская сетка, 74
- чебышёвские узлы, 74
- численное дифференцирование, 99
- число обусловленности, 263
- дефект сплайна, 86
- диагонализуемая матрица, 390
- диагональное преобладание, 222
- дифференциальная арифметика, 118
- дифференцирование
 - алгоритмическое, 99, 118
 - дифференцирование
 - автоматическое, 99, 118
 - дифференцирование численное, 99
 - дифференцирование символьное, 99
 - длина вектора, 226
 - доминирующее собственное значение, 404
 - доминирующий собственный вектор, 404
 - естественный сплайн, 95
 - евклидова норма, 226
 - экспоненциальная трудоёмкость, 34
- экстраполяция, 67
- экстремум глобальный, 359
- экстремум локальный, 359
- эквивалентные нормы, 230, 243
- элементарная матрица
 - перестановок, 284
- энергетическая норма, 245
- энергии функционал, 357
- эрмитова интерполяция, 76
- формула Ньютона-Лейбница, 142
- формула Родрига, 135
- формула Симпсона, 152
- формула кубатурная, 144
- формула квадратурная, 144
- формула парабол, 152
- формула прямоугольников, 146
- формула трапеций, 149

- формулы Гаусса, 166
формулы Лобатто, 180
формулы Маркова, 180
формулы Ньютона-Котеса, 146
формулы численного
дифференцирования,
101, 103
функционал энергии, 357
главный элемент, 283
гёльдерова норма, 226
характеристическое уравнение
матрицы, 208
характеризация Бека, 470
характеризация Оетгли-Прагера,
470
хессенбергова форма, 402
индуцированная норма, 240
интегральная метрика, 43
интерполирование, 44
интерполяционная квадратурная
формула, 157
интерполяция эрмитова, 76
интерполянт, 44
интервальная арифметика, 23
интервальное расширение, 27
итерационные методы, 274
каноническая форма СЛАУ, 272
каноническая форма Самарского,
374
классическая интервальная
арифметика, 24
коэффициент чувствительности,
20
коэффициенты Фурье, 127
коэффициенты перекоса, 394
коллинеарные векторы, 202
комплексификация, 249
конечные методы, 274
кратность узла, 75
круги Гершгорина, 398
квадратурная интерполяционная
формула, 157
линейная интерполяция, 49
линейная оболочка, 202
линейная задача о наименьших
квадратах, 385
максимум-норма, 226
машинная интервальная
арифметика, 37
матрица Гильберта, 131, 267
матрица Грама, 127
матрица Уилкинсона, 396
матрица Вандермонда, 49, 268
матрица диагонализуемая, 390
матрица наклонов интервальная,
478
матрица недефектная, 390
матрица неособенная, 204
матрица неразложимая, 224
матрица особенная, 204
матрица отражения, 303
матрица перестановок, 286
матрица почти треугольная, 402
матрица предобуславливающая,
334
матрица простой структуры, 390
матрица разложимая, 224
матрица регулярная, 204
матрица скалярная, 336
матрица строго нижняя
треугольная, 341
матрица строго регулярная, 288
матрица строго верхняя
треугольная, 341
матрица транспозиции, 285
матрица трёхдиагональная, 318
матрица вращения, 310
матричная норма, 235
матричный ряд Неймана, 253
мера диагонального
преобладания, 348
метод Эйлера, 372
метод Гаусса, 278
метод Гаусса-Зейделя, 345

- метод Герона, 463
метод Хаусхолдера, 306
метод Холесского, 295
метод Шульца, 380
метод Якоби, 340
метод градиентного спуска, 359
метод квадратного корня, 296
метод минимальных невязок, 367
метод наискорейшего спуска, 363
метод отражений, 306
метод прогонки, 320
метод простой итерации, 336
метод релаксации, 351
метод сопряжённых градиентов, 371
метод установления, 372
метод ветвлений и отсечений, 483
метрика, 43, 119
множитель Холесского, 291
мультиметрика, 461
насыщение численного метода, 95
натуральный сплайн, 95
недефектная матрица, 390
нелинейная интерполяция, 49
ненасыщаемый метод, 95, 180
непрерывность по Липшицу, 21, 84
неравенство Коши-Буняковского, 226
неравенство Минковского, 226
нестационарный итерационный процесс, 324
невязка, 325, 350
норма, 225
норма энергетическая, 245
норма индуцированная, 240
норма операторная, 240
норма подчинённая, 240
нормальная система уравнений, 385
обобщённая степень, 62
обратные степенные итерации, 413
оператор Кравчика, 480
оператор Ньютона интервальный, 475
операторная форма СЛАУ, 274
операторная норма, 240
ортогонализация Грама-Шмидта, 133, 313
основная теорема интервальной арифметики, 28
остаточный член квадратурной формулы, 144
относительная погрешность, 12
отношение Рэлея, 399
почти решения, 441
подчинённая норма, 240
подпространства Крылова, 316
погрешность абсолютная, 11
погрешность относительная, 12
поле значений матрицы, 400
полином интерполяционный, 48
полином интерполяционный Лагранжа, 52
полином интерполяционный Ньютона, 61
полиномы Чебышёва, 68
полиномы Лежандра, 134
полиномиальная трудоёмкость, 34
порядок аппроксимации, 106
порядок точности формулы, 106, 181
правило Рунге, 193
предобуславливание, 334
предобуславливатель, 334
пример Бернштейна, 82
пример Рунге, 82
принцип релаксации, 349
принцип вариационный, 355
приведённые полиномы Чебышёва, 72
признак Адамара, 223
пространство строго

- нормированное, 123
- прямые методы, 274
- псевдометрика, 44
- псевдорасстояние, 44
- расстояние, 43, 119
- расщепление матрицы, 335
- равномерная метрика, 43
- разделённая разность, 52
- разложение Холесского, 291
- разложение Шура, 213
- разложение сингулярное, 220
- разностные уравнения
 - трёхточечные, 319
- разность назад, 101
- разность вперёд, 101
- рекуррентный вид системы, 326, 434
- рекуррентный вид уравнения, 434
- ряд Фурье, 127
- сдвиг спектра, 414
- сетка, 44, 144
- схема единственного деления, 279
- сходимость по норме, 230, 243
- сходимость поэлементная, 244
- символьное дифференцирование, 99
- сингулярные числа, 214
- сингулярные векторы, 214
- система трёхдиагональная, 318
- скалярные произведения, 202
- след матрицы, 419
- собственный вектор, 385
- собственное значение, 385
- спектр матрицы, 208
- спектральная норма, 242
- спектральный радиус, 247
- сплайн, 86
- среднеквадратичная метрика, 43
- стационарный итерационный процесс, 324
- степенной метод, 407
- степень сплайна, 86
- строго нормированное пространство, 123
- строго регулярная матрица, 288
- субдистрибутивность, 25
- сжатие, 460
- сжимающее отображение, 460
- шаблон, 103
- теорема Абеля-Рурфини, 388
- теорема Банаха о неподвижной точке, 461
- теорема Бауэра-Файка, 390
- теорема Больцано-Коши, 456
- теорема Брауэра о неподвижной точке, 471
- теорема Экарта-Янга, 259
- теорема Фабера, 84
- теорема Гершгорина, 398
- теорема Леви-Деспланка, 224
- теорема Марцинкевича, 85
- теорема Миранды, 457
- теорема Островского, 388
- теорема Островского-Райха, 354
- теорема Самарского, 376
- теорема Стеклова-Пойа, 185
- теорема Шрёдера о неподвижной точке, 462
- теорема Таусски, 224
- теорема Вейерштрасса, 81
- теорема Вейля, 401
- теорема Виландта-Хофмана, 427
- теорема о сингулярном разложении, 220
- тест существования решения, 474
- треугольное разложение, 282
- тригонометрические полиномы, 47
- трёхдиагональная матрица, 94
- узлы сплайна, 86
- ведущая подматрица, 203
- ведущий элемент, 282
- ведущий минор, 203
- векторная норма, 225

- вырожденный интервал, 23
- задача некорректная, 19, 117
- задача о наименьших квадратах
 - линейная, 385
- задача приближения функции,
 - 121
- задача сглаживания, 121
- задача вычислительно
 - корректная, 437
- значащая цифра, 12
- P -сжатие, 461
- 1-норма, 226

- LDL-разложение, 297
- LU-разложение, 282

- О-большое, 95

- QR-алгоритм, 423, 426
- QR-разложение, 300