



А. А. Амосов, Ю. А. Дубинский, Н. В. Копченова

ВЫЧИСЛИТЕЛЬНЫЕ МЕТОДЫ



$$\begin{matrix} a_{11} & a_{12} & 0 & 0 & \dots & 0 \\ a_{21} & a_{22} & a_{23} & 0 & \dots & 0 \\ 0 & a_{32} & a_{33} & a_{34} & \dots & 0 \end{matrix}$$

$$\begin{matrix} 0 & 0 & \dots & 0 & a_{m-1, m-2} & a_{m-1, m-1} & a_{m-1, m} \\ 0 & 0 & \dots & 0 & \dots & 0 & a_{m, m-1} \end{matrix}$$

$$\begin{matrix} 0 & 0 & \dots & 0 & a_{m, m-1} & a_{m, m} \end{matrix}$$

**А. А. АМОСОВ,
Ю. А. ДУБИНСКИЙ,
Н. В. КОПЧЕНОВА**

ВЫЧИСЛИТЕЛЬНЫЕ МЕТОДЫ

Учебное пособие
Издание четвертое, стереотипное



САНКТ-ПЕТЕРБУРГ • МОСКВА • КРАСНОДАР
2014

ББК 22.1я73

А 62

Амосов А. А., Дубинский Ю. А., Копченова Н. В.

А 62 Вычислительные методы: Учебное пособие. — 4-е изд., стер. — СПб.: Издательство «Лань», 2014. — 672 с.: ил. — (Учебники для вузов. Специальная литература).

ISBN 978-5-8114-1623-3

В книге рассматриваются вычислительные методы, наиболее часто используемые в практике прикладных и научно-технических расчетов: методы решения задач линейной алгебры, нелинейных уравнений, проблемы собственных значений, методы теории приближения функций, численное дифференцирование и интегрирование, поиск экстремумов функций, решение обыкновенных дифференциальных уравнений, численное решение интегральных уравнений, линейная и нелинейная задачи метода наименьших квадратов, метод сопряженных градиентов. Значительное внимание уделяется особенностям реализации вычислительных алгоритмов на компьютере и оценке достоверности полученных результатов. Имеется большое количество примеров и геометрических иллюстраций. Даются сведения о стандарте IEEE, о сингулярном разложении матрицы и его применении для решения переопределенных систем, о двухслойных итерационных методах, о квадратурных формулах Гаусса–Кронрода, о методах Рунге–Кутты–Фельберга.

Учебное пособие предназначено для студентов всех направлений подготовки, обучающихся в классических и технических университетах и изучающих вычислительные методы, будет полезно аспирантам, инженерам и научным работникам, применяющим вычислительные методы в своих исследованиях.

ББК 22.1я73

Рецензенты:

Н. Н. КАЛИТКИН — член-корреспондент РАН;

Н. С. БАХВАЛОВ — профессор, зав. кафедрой вычислительной математики механико-математического факультета МГУ, академик РАН;

Е. И. МОИСЕЕВ — профессор, декан факультета вычислительной математики и кибернетики МГУ, член-корреспондент РАН.

**Обложка
Е. А. ВЛАСОВА**

© Издательство «Лань», 2014
© Коллектив авторов, 2014
© Издательство «Лань»,
художественное оформление, 2014

Цель расчетов — не числа, а понимание.

P.B. Хемминг

Из нашего девиза «Цель расчетов — не числа, а понимание» следует, что человек, который должен этого понимания достигнуть, обязан знать, как происходит вычисление. Если он не понимает, что делается, то очень маловероятно, чтобы он извлек из вычислений что-нибудь ценное. Он видит голые цифры, но их истинное значение может оказаться скрытым в вычислениях.

P.B. Хемминг

ПРЕДИСЛОВИЕ К ПЕРВОМУ ИЗДАНИЮ

В настоящее время имеется значительное число учебников и монографий, посвященных методам вычислений (часть из них отражена в списке литературы к данному пособию). Однако, на наш взгляд, большинство этих книг ориентировано на студентов-математиков или на специалистов по вычислительной математике. В то же время практически отсутствует отечественная учебная литература, в которой доступным для студента технического вуза или инженера образом были бы изложены основы вычислительных методов, применяемых сегодня для решения инженерных задач. Особенно острой, по мнению авторов, является потребность в книге, которая содержала бы не только изложение начал численных методов, но и давала бы представление о реально используемых в вычислительной практике алгоритмах. Данное учебное пособие призвано в определенной степени восполнить этот пробел.

Настоящее пособие адресовано в первую очередь студентам и аспирантам высших технических учебных заведений, изучающим основы математического моделирования и численные методы. Авторы надеются на то, что эта книга будет полезна широкому кругу инженерных и научно-технических работников, которые намерены применять ЭВМ для решения прикладных задач.

При написании книги авторы использовали многолетний опыт преподавания курса вычислительных методов студентам и аспирантам различных специальностей Московского энергетического института, а также опыт работы в вычислительном центре МЭИ. Значительное влияние на выбор материала и характер изложения оказали также многочисленные дискуссии со слушателями существующего при МЭИ факультета повышения квалификации преподавателей вузов страны.

Авторы стремились изложить материал по возможности наиболее простым и доступным образом. Объем знаний высшей математики, необходимый для понимания содержания книги, не выходит за рамки программы младших курсов втуза. Пособие содержит довольно много

примеров, иллюстрирующих те или иные положения теории, а также демонстрирующих особенности вычислительных методов или работу конкретных алгоритмов. Тем не менее, многие из рассматриваемых вопросов трудны для восприятия и требуют внимательного изучения. К сожалению, в учебной литературе они нередко опускаются. К таким центральным вопросам относятся, например понятия корректности, устойчивости и обусловленности вычислительных задач и вычислительных алгоритмов, особенности поведения вычислительной погрешности. Важность их понимания для эффективного применения ЭВМ сегодня велика и не акцентировать на них внимание авторы посчитали невозможным.

Дадим краткое изложение основного содержания книги. Важную идеиную нагрузку несут на себе первые три главы. Рассматриваемые в них вопросы закладывают фундамент, необходимый для правильного понимания изложенных в остальных главах вычислительных методов.

В гл. 1 дается общее представление о методе математического моделирования, в том числе о процессе создания математических моделей, последовательности этапов решения инженерной задачи с применением ЭВМ, вычислительном эксперименте.

В гл. 2 наряду с введением в элементарную теорию погрешностей содержится изложение основных особенностей машинной арифметики. Понимание этих особенностей необходимо тем, кто заинтересован в эффективном применении ЭВМ для решения прикладных задач.

В гл. 3 обсуждаются важнейшие свойства вычислительных задач, методов и алгоритмов. Дается общее представление о корректности, устойчивости и обусловленности вычислительной задачи. Приводится описание основных классов вычислительных методов. Значительное внимание уделяется устойчивости вычислительных алгоритмов, их чувствительности к ошибкам. Дается представление о различных подходах к анализу ошибок, в том числе и об обратном анализе ошибок. Обсуждаются требования, предъявляемые к вычислительным алгоритмам.

Конкретные вычислительные задачи и методы их решения рассматриваются начиная с гл. 4. Здесь авторы стремились к тому, чтобы не только изложить простейшие подходы, но и дать представление об алгоритмах, которые реально используются для решения соответствующих задач.

В гл. 4 рассматриваются методы отыскания решений нелинейных уравнений. Значительное вниманиеделено постановке задачи и ее свойствам, в частности — чувствительности корней нелинейных уравнений к погрешностям. Среди различных методов отыскания корней более подробно излагаются метод простой итерации, метод Ньютона и различные их модификации.

В гл. 5 рассмотрены прямые (точные) методы решения систем линейных алгебраических уравнений. Основное внимание уделяется

методу Гаусса и его различным модификациям. Рассматриваются использование *LU*-разложения матриц для решения систем линейных уравнений, метод квадратных корней, метод прогонки, методы вращений и отражений. Обсуждается алгоритм итерационного уточнения.

В гл. 6 рассматриваются итерационные методы решения систем линейных алгебраических уравнений: метод простой итерации, метод Зейделя, метод релаксации и другие методы.

В гл. 7 рассматривается задача отыскания решений систем нелинейных уравнений. Обсуждаются не только соответствующие итерационные методы, но и различные подходы к решению сложной задачи локализации.

В гл. 8 дается представление о проблеме собственных значений и о различных подходах к вычислению собственных значений и собственных векторов. Излагаются степенной метод и обратный степенной метод, обсуждается *QR*-алгоритм.

В гл. 9 излагаются наиболее известные численные методы решения задачи одномерной минимизации, в том числе метод деления отрезка пополам, метод Фибоначчи, метод золотого сечения и метод Ньютона.

В гл. 10 рассматриваются различные методы решения задачи безусловной минимизации. Наиболее полно изложены градиентный метод, метод Ньютона и метод сопряженных градиентов.

В гл. 11 рассмотрены наиболее важные и часто встречающиеся в приложениях методы приближения функций. Значительное внимание удалено интерполяции, причем рассматривается интерполяция не только алгебраическими многочленами, но и тригонометрическими многочленами, а также интерполяция сплайнами. Достаточно подробно обсуждается метод наименьших квадратов. Даётся понятие о наилучшем равномерном приближении и дробно-рациональных аппроксимациях.

В эту главу включены также некоторые вопросы, имеющие непосредственное отношение к методам приближения функций. Это конечные и разделенные разности, многочлены Чебышева, быстрое дискретное преобразование Фурье.

В гл. 12 рассматриваются различные подходы к выводу формул численного дифференцирования, обсуждается чувствительность этих формул к ошибкам в вычислении значений функции.

В гл. 13 излагаются методы вычисления определенных интегралов. Выводятся квадратурные формулы интерполяционного типа и квадратурные формулы Гаусса. Даётся представление о принципах построения адаптивных процедур численного интегрирования и, в частности, об используемых в них способах апостериорной оценки погрешности. Рассматриваются различные подходы к вычислению интегралов от функций, имеющих те или иные особенности. В частности, затрагивается проблема интегрирования быстро осциллирующих функций.

Глава 14 посвящена численным методам решения задачи Коши для обыкновенных дифференциальных уравнений. Подробно рассматриваются метод Эйлера и его различные модификации. Значительное внимание уделено рассмотрению классических методов Рунге—Кутты и Адамса. Обсуждаются различные свойства устойчивости численных методов решения задачи Коши, в том числе нуль-устойчивость, абсолютная устойчивость, A -устойчивость, $A(\alpha)$ -устойчивость. Специально рассматриваются жесткие задачи и методы их решения.

В гл. 15 изучаются методы численного решения двухточечных краевых задач. Подробно излагается применение метода конечных разностей к решению краевых задач для обыкновенного дифференциального уравнения второго порядка. Дается представление о проекционных методах Ритца и Галеркина; обсуждается один из их современных вариантов — метод конечных элементов. Завершает главу рассмотрение метода пристрелки.

Можно предположить, что к изучению отдельных параграфов или даже глав этой книги читатель приступит, только столкнувшись с необходимостью решить на ЭВМ важную для него задачу. Вероятно, в этом случае польза от изучения соответствующего раздела будет наибольшей. Многие вопросы рассмотрены очень кратко или не рассмотрены вообще. Авторы надеются на то, что соответствующие пробелы можно восполнить, использовав сделанные в тексте ссылки на известные учебники и монографии.

В ряде случаев авторы позволяют себе интерпретировать те или иные результаты, делать определенные выводы и даже давать рекомендации в надежде на то, что для новичка соответствующие рассуждения дадут полезный начальный ориентир. В целом же ко всяkim рекомендациям в такой сложной и многообразной области, как применение вычислительных методов для решения прикладных задач на ЭВМ, следует отнестись с осторожностью. Они не могут претендовать на бесспорность и их следует рассматривать скорее как отражение точки зрения авторов.

Иногда то или иное положение обосновывается ссылкой на вычислительную практику. Хотя критерий практики и играет при отборе методов вычислений существенную роль, все же оценки методов, основанные на результатах их применения для решения конкретных задач, нередко бывают весьма субъективны и противоречивы.

В заключение отметим, что никакие теоретические положения и советы не могут заменить собственного опыта вычислительной работы. Как надеются авторы, параллельно с изучением данной книги такой опыт может приобрести читатель, переходя от решения задач учебного характера к серьезным практическим задачам.

Авторы

ПРЕДИСЛОВИЕ КО ВТОРОМУ ИЗДАНИЮ

Со времени выхода в свет первого издания этой книги, в 1994 г. в издательстве «Высшая школа», прошло около десяти лет. За это время стало ясно, что она нашла своего читателя. Книга активно используется студентами, аспирантами, преподавателями и научными сотрудниками не только Московского энергетического института, но и других вузов страны. В то же время обнаружился ряд опечаток и неточностей в изложении, которые было бы желательно исправить. Поэтому авторы с благодарностью приняли любезное предложение Издательства МЭИ о переиздании учебного пособия.

Во втором издании книги исправлены замеченные недостатки. Рассмотрен ряд вопросов, которые не были включены в первое издание. Это стандарт IEEE, сингулярное разложение матрицы и его применение для решения переопределенных систем, двухслойные итерационные методы решения систем линейных алгебраических уравнений (включая метод с чебышевским набором параметров, метод сопряженных градиентов и предобусловливание), обратная интерполяция, интерполяционный многочлен Бесселя, многомерная интерполяция, квадратурные формулы Гаусса—Кронрода, метод Рунге—Кутты—Фельберга. Дополнен список литературы.

Авторы пользуются случаем выразить свою благодарность российским ученым, которые оказали существенное влияние на формирование их взглядов на вычислительную математику. Это акад. РАН Н.С. Бахвалов, акад. РАН В.В. Воеводин, акад. РАН А.А. Самарский, акад. РАН А.Н. Тихонов, проф. Ф.П. Васильев, проф. Е.Г. Дьяконов, проф. Я.М. Жилийкин, проф. Х.Д. Икрамов, проф. Г.М. Кобельков, докт физ-мат. наук А.С. Кронрод, проф. И.А. Марон.

Работа по подготовке второго издания к печати выполнена проф. А.А. Амосовым. Авторы искренне признательны своим коллегам проф. А.А. Злотнику, доц. О.А. Амосовой и доц В.П. Григорьеву, чьи замечания были существенно использованы в процессе работы над вторым изданием. Авторы заранее благодарны всем, кто пожелает высказать свои замечания и предложения. Их следует направлять по адресу: 111250, Москва, Красноказарменная ул., д. 14, Издательство МЭИ или по адресу AmosovAA@mpei.ru.

Авторы

ПРЕДИСЛОВИЕ К ТРЕТЬЕМУ ИЗДАНИЮ

Перед Вами третье издание книги, которая теперь называется «Вычислительные методы», хотя в двух предыдущих изданиях она выходила под названием «Вычислительные методы для инженеров». По-видимому, авторы должны как-то объяснить, почему название книги стало существенно короче.

Во-первых, методы вычислений, применяемые для решения возникающих на практике задач, универсальны. И не существует каких-то специальных методов, предназначенных для решения только лишь инженерных задач. Поэтому изложенные в книге методы могут использоваться студентами, аспирантами и исследователями для решения самых разных задач, в том числе — экономических, медицинских, биологических и т.д. Нет необходимости специально сужать круг возможных читателей.

Во-вторых, перемены, происходящие в высшем образовании нашей страны, скоро приведут к тому, что высокое звание «инженер» станет восприниматься как архаика. Книга, предназначенная для каких-то непонятных инженеров, может показаться ненужной будущим бакалаврам и магистрам. Так что сокращение названия книги — это и вынужденная дань нынешней моде.

Чем это издание отличается от предыдущего? Во-первых, в нем исправлены обнаруженные неточности и опечатки. Во-вторых, добавлена новая глава, посвященная численному решению интегральных уравнений. В третьих, частично изменено содержание некоторых параграфов. Существенной переработке подвергся параграф 2.5, посвященный особенностям машинной арифметики. В нем теперь достаточно полно отражены особенности IEEE арифметики, реализованной в настоящее время на большинстве современных компьютеров. Достойное место в книге, наконец, занял метод наименьших квадратов. В ней появились два новых параграфа, посвященных линейной и нелинейной задачам метода наименьших квадратов. Появился также отдельный параграф, посвященный методу сопряженных градиентов решения систем линейных алгебраических уравнений.

Список литературы дополнен. Ссылки на литературу, добавленную к третьему изданию, помечены звездочкой.

Работа по подготовке к печати третьего издания выполнена проф. А.А. Амосовым.

Авторы заранее благодарны читателям, которые пожелаю высказать свои замечания или предложения по содержанию книги. Их следует направлять по адресу 111250, Москва, Красноказарменная ул., д. 14, Издательский дом МЭИ или по адресу электронной почты AmosovAA@trei.ru.

Авторы

Глава 1

МАТЕМАТИЧЕСКОЕ МОДЕЛИРОВАНИЕ И РЕШЕНИЕ ПРИКЛАДНЫХ ЗАДАЧ С ПРИМЕНЕНИЕМ КОМПЬЮТЕРА

В этой главе дается общее представление о методе математического моделирования, в том числе о процессе создания математических моделей, о последовательности этапов решения прикладной задачи с применением компьютера, о вычислительном эксперименте.

Разделение решаемых с применением компьютера задач на прикладные и научные является до известной степени условным. Тем не менее, при написании данного пособия авторы ориентировались на читателя, интересующегося решением именно прикладных задач. Попытаемся охарактеризовать этот класс задач, выделив его некоторые характерные особенности.

1. Прикладные задачи имеют ярко выраженную практическую направленность. Целью их решения является создание новой конструкции, разработка нового технологического процесса, минимизация затрат на производство некоторого изделия и т.д. Поэтому для таких задач характерна необходимость доведения результатов до конкретных чисел, графиков, таблиц, на основании которых можно принимать решения.

2. Эти задачи характеризуются значительным объемом выполняемой вычислительной работы.

3. Для этих задач характерно использование достаточно сложных математических моделей и серьезного математического аппарата.

4. Прикладные задачи, как правило, решают специалисты, не являющиеся профессионалами в области разработки математических методов и программного обеспечения компьютеров. Поэтому естественно желание этих специалистов использовать готовые вычислительные методы и стандартное математическое программное обеспечение.

Наконец, условимся считать, что в рассматриваемый класс задач входят задачи только умеренной сложности. Для их решения не требуется сверхбольшие скорости вычислений и сверхбольшие объемы памяти для хранения данных. Таким образом, эти задачи могут быть решены с помощью компьютеров, доступных массовому пользователю. Те же задачи, которые требуют для решения сверхмощной вычислительной техники и принципиально новых алгоритмов, будем относить к категории научных задач.

§ 1.1. Математическое моделирование и процесс создания математической модели

Математическое моделирование представляет собой метод исследования объектов и процессов реального мира с помощью их приближенных описаний на языке математики — *математических моделей*. Этот метод чрезвычайно плодотворен и известен уже несколько тысячелетий. Насущные задачи земледелия и строительства еще в древние времена приводили к необходимости определения площадей и объемов, а следовательно, и к рассмотрению элементарных геометрических фигур, дающих пример простейших математических моделей. Возможности математического моделирования и его влияния на научно-технический прогресс неизмеримо возросли в последние десятилетия в связи с созданием и широким внедрением компьютеров.

Процесс создания математической модели условно можно разбить на ряд основных этапов:

- 1) построение математической модели;
- 2) постановка, исследование и решение соответствующих вычислительных задач,
- 3) проверка качества модели на практике и модификация модели.

Рассмотрим основное содержание этих этапов.

1. Построение математической модели. Предполагается, что с помощью наблюдений и экспериментов, практики (понимаемой в самом широком смысле) получена достаточно подробная информация об изучаемом явлении. Для рассматриваемого этапа характерно глубокое проникновение в полученные факты в целях выяснения главных закономерностей. Выявляются основные «характеристики» явления, которым сопоставляются некоторые величины. Как правило, эти величины принимают числовые значения, т.е. являются переменными, векторами, матрицами, функциями и т.д.

Установленным внутренним связям между «характеристиками» явления придается форма равенств, неравенств, уравнений и логических структур, связывающих величины, включенные в математическую модель. Таким образом, математическая модель становится записью на языке математики законов природы, управляющих протеканием исследуемого процесса или описывающих функционирование изучаемого объекта. Она включает в себя набор некоторых величин и описание характера связи между ними.

Построение математических моделей — существенная и очень важная часть естественных и технических наук. Эта задача, требующая от иссле-

дователя глубокого знания предметной области, высокой математической культуры, опыта построения моделей, развитой интуиции и многое другое. Создание удачной новой модели — всегда крупное достижение соответствующей науки, а иногда и целый этап в ее развитии.

Подчеркнем, что математическая модель неизбежно представляет собой компромисс между бесконечной сложностью изучаемого явления и желаемой простотой его описания. Модель должна быть достаточно полной, для того чтобы оказаться полезной для изучения свойств исследуемого явления. В то же время она обязана быть достаточно простой, для того чтобы допускать возможность ее анализа существующими в математике средствами и ее реализации на компьютере. Из огромного числа характеристик явления и действующих на него факторов требуется выделить основные, определяющие, отбросив при этом второстепенные, несущественные.

Нередко в математическую модель закладываются некоторые гипотезы, еще не подтвержденные на практике. Такую математическую модель часто называют *гипотетической*.

Приведем пример простейшей математической модели.

Пример 1.1. Пусть исследуется движение тела, брошенного со скоростью v_0 под углом α к поверхности Земли.

Будем считать, что в рассматриваемом случае можно пренебречь сопротивлением воздуха, считать Землю плоской, а ускорение свободного падения g — постоянной. Введем систему координат, ее начало поместим в точку бросания, ось Ox направим горизонтально в направлении бросания, а ось Oy — вертикально вверх (рис. 1.1). Пусть $u(t)$ и $w(t)$ — горизонтальная и вертикальная составляющие скорости $v(t)$ в момент времени t (в начальный момент $t = 0$, $v = v_0$)

Согласно законам механики, при сделанных предположениях движение тела в горизонтальном направлении является равномерным, а в вер-

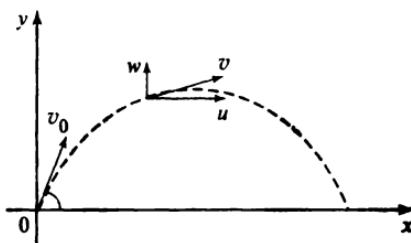


Рис. 1.1

тикальном — равноускоренным с ускорением, равным $-g$. Поэтому справедливы следующие равенства:

$$u = v_0 \cos\alpha, \quad x = (v_0 \cos\alpha)t, \quad (1.1)$$

$$w = v_0 \sin\alpha - gt, \quad y = (v_0 \sin\alpha)t - \frac{gt^2}{2}. \quad (1.2)$$

Формулы (1.1), (1.2) и дают простейшую математическую модель рассматриваемого явления, созданную в XVII в. Г. Галилеем¹. Заметим, что при $0 < \alpha < \pi/2$ траектория движения представляет собой параболу

$$y = -\frac{g}{2v_0^2 \cos^2\alpha} x^2 + (\tan\alpha)x. \quad \blacktriangle$$

Математические модели часто разделяют на статические и динамические. Статическая модель описывает явление или ситуацию в предположении их завершенности, неизменности (т.е. в статике). Динамическая модель описывает, как протекает явление или изменяется ситуация от одного состояния к другому (т.е. в динамике). При использовании динамических моделей, как правило, задают начальное состояние системы, а затем исследуют изменение этого состояния во времени.

2. Постановка, исследование и решение вычислительных задач. Для того чтобы найти интересующие исследователя значения величин или выяснить характер из зависимости от других входящих в математическую модель величин, ставят, а затем решают математические задачи.

Выявим основные типы решаемых задач. Для этого все величины, включенные в математическую модель, условно разобьем на три группы:

- 1) исходные (входные) данные x ;
- 2) параметры модели a ;
- 3) искомое решение (выходные данные) y .

В динамических моделях искомое решение часто является функцией времени $y = y(t)$; переменная t в таких моделях, как правило, бывает выделенной и играет особую роль.

Наиболее часто решают так называемые *прямые задачи*, постановка которых выглядит следующим образом: по данному значению входного данного x при фиксированных значениях параметров a требуется найти решение y . Процесс решения прямой задачи можно рассматривать как математическое моделирование причинно-следственной связи, присущей явлению. Тогда входное данное x характеризует «причины» явле-

¹ Галилео Галилей (1564—1642) — итальянский физик, механик, астроном, один из основателей точного естествознания.

ния, которые задаются и варьируются в процессе исследования, а искомое решение y — «следствие».

Для того чтобы математическое описание было применимо не к единичному явлению, а к широкому кругу близких по природе явлений, в действительности строят не единичную математическую модель, а некоторое параметрическое семейство моделей. Будем считать, что выбор конкретной модели из этого семейства осуществляется фиксацией значений параметров модели a . Например, в роли таких параметров могут выступать некоторые из коэффициентов, входящих в уравнения. С помощью выбора параметров может производиться указание типа функциональной зависимости между некоторыми из величин. На конец, если используемые математические модели разбиты на классы, то параметром может служить и класс используемой модели.

Пример 1.2. Для модели (1.1), (1.2) прямую задачу естественно формулировать как задачу вычисления величин $u(t)$, $w(t)$, $x(t)$, $y(t)$ по задаваемым входным данным v_0 , α . Параметром модели здесь является ускорение свободного падения g . Его значение зависит от того, производится ли бросание тела с поверхности Земли на уровне Мирового океана, в глубокой шахте или же на большой высоте. Заметим, что та же модель пригодна для описания движения тела, брошенного на любой другой планете, если значение параметра g для этой планеты известно. ▲

Большую роль играет решение так называемых *обратных задач*, состоящих в определении входного x по данному значению y (параметры модели a , как и в прямой задаче, фиксированы). Решение обратной задачи — это в определенном смысле попытка выяснить, какие «причины» x привели к известному «следствию» y . Как правило, обратные задачи оказываются сложнее для решения, чем прямые.

Пример 1.3. Для модели (1.1), (1.2) обратную задачу можно сформулировать так: по заданным $u(t)$, $w(t)$, $x(t)$, $y(t)$ требуется найти значения v_0 , α . Заметим, что для однозначного определения v_0 , α достаточно задать в любой фиксированный момент $t_0 \geq 0$ одну из пар величин $(u(t_0), w(t_0))$ или $(x(t_0), y(t_0))$. ▲

Помимо двух рассмотренных типов задач следует упомянуть еще один тип — *задачи идентификации*. В широком смысле задача идентификации модели — это задача выбора среди множества всевозможных моделей той, которая наилучшим образом описывает изучаемое явление. В такой постановке эта задача выглядит как практически неразрешимая проблема. Чаще задачу идентификации понимают в узком смысле, как задачу выбора из заданного параметрического семейства моделей

конкретной математической модели (с помощью выбора ее параметров a), с тем чтобы оптимальным в смысле некоторого критерия образом согласовать следствия из модели с результатами наблюдений.

Пример 1.4. Применительно к модели (1.1), (1.2) задача идентификации может состоять в определении значения ускорения свободного падения планеты g по результатам наблюдений за параметрами траектории. ▲

Указанные три типа задач (прямые, обратные и задачи идентификации) будем называть *вычислительными задачами*. Для удобства изложения в дальнейшем независимо от типа решаемой задачи будем называть набор подлежащих определению величин *искомым решением* и обозначать через y , а набор заданных величин — *входным данным* и обозначать через x .

Пример 1.5. При описании многих явлений используют модель полиномиальной зависимости между величинами x и y :

$$y = P_n(x) \equiv a_0 + a_1 x + \dots + a_n x^n. \quad (1.3)$$

Здесь a_0, a_1, \dots, a_n — коэффициенты многочлена, являющиеся параметрами модели (в число параметров модели можно включить и степень многочлена).

При фиксированных значениях параметров прямая задача состоит в вычислении значения многочлена $y = P_n(x)$ по заданному x . В таком случае целью решения обратной задачи является определение по заданному значению y соответствующего ему значения x . Нетрудно видеть, что это есть задача отыскания корней многочлена, отличающегося от $P_n(x)$ заменой коэффициента a_0 на $\tilde{a}_0 = a_0 - y$. Если же из практики известна некоторая информация о зависимости y от x , то определение параметров a_0, a_1, \dots, a_n , при которых модель (1.3) наилучшим в некотором смысле образом описывает эту зависимость, представляет собой задачу идентификации. Например, если задана таблица значений x_i, y_i , ($i = 1, \dots, N$), то такую задачу в зависимости от ситуации можно решать, используя известные методы интерполяции и наименьших квадратов (см. гл. 11). ▲

Пример 1.6. Нередко входящие в модель функции $x(t)$ и $y(t)$ бывают связаны равенством

$$y(t) = y_0 + \int_0^t x(\tau) d\tau.$$

Например, так связаны между собой скорость $x(t)$ и путь $y(t)$ при прямолинейном движении. Тогда при фиксированном значении посто-

янной y_0 прямая задача (задача интегрирования) состоит в вычислении первообразной $y(t)$ по заданной функции $x(t)$. Обратная задача (задача дифференцирования) заключается в вычислении $x(t) = y'(t)$ по заданной функции $y(t)$. ▲

Как правило, решение вычислительной задачи не удается выразить через входные данные в виде конечной формулы. Однако это совсем не означает, что решение такой задачи не может быть найдено. Существуют специальные методы, которые называют *численными* (или *вычислительными*). Они позволяют свести получение численного значения решения к последовательности арифметических операций над численными значениями входных данных. Эти методы были известны давно: в качестве примера, уже ставшего классическим, можно привести открытие Леверье¹ в 1846 г. новой планеты Нептун. Однако для решения задач численные методы применялись довольно редко, так как их использование предполагает выполнение гигантского объема вычислений². Поэтому в большинстве случаев до появления компьютеров приходилось избегать использования сложных математических моделей и исследовать явления в простейших ситуациях, когда возможно найти аналитическое решение. Несовершенство вычислительного аппарата становилось фактором, сдерживающим широкое использование математических моделей в науке и технике.

Появление компьютеров кардинально изменило ситуацию. Класс математических моделей, допускающих подробное исследование, резко расширился. Решение многих, еще недавно недоступных, вычислительных задач стало обыденной реальностью.

3. Проверка качества модели на практике и модификация модели. На этом этапе выясняют пригодность математической модели для описания исследуемого явления. Теоретические выводы и конкретные результаты, вытекающие из гипотетической математической модели, сопоставляют с экспериментальными данными. Если они противоречат друг другу, то выбранная модель непригодна и ее следует пересмотреть, вернувшись к первому этапу. Если же результаты совпадают с допустимой для описания данного явления точностью, то модель можно при-

¹ Урбен Жан Жозеф Леверье (1811—1877) — французский астроном. На основании законов небесной механики, использовав данные об аномалиях в движении планеты Уран, Леверье рассчитал траекторию движения гипотетической неизвестной планеты. В том же году немецкий астроном Галле обнаружил Нептун в указанном Леверье месте.

² Расчет траектории планеты Нептун потребовал от Леверье нескольких месяцев кропотливой вычислительной работы.

знать пригодной. Конечно, необходимо дополнительное исследование в целях установления степени достоверности модели и границ ее применимости.

На определенном этапе развития науки и техники постепенное накопление знаний приводит к моменту, когда результаты, получаемые с помощью математической модели, вступают в противоречие с данными практики или перестают удовлетворять ее требованиям в смысле точности. Тогда возникает необходимость модификации модели или же создания принципиально новой, более сложной модели. Таким образом, цикл создания математической модели повторяется многократно.

Пример 1.7. Рассмотрим задачу внешней баллистики, т.е. задачу о движении артиллерийского снаряда. Простейшая модель (1.1), (1.2) дает параболическую траекторию движения снаряда, что, как было замечено еще в XVII в., противоречит данным практики. Существенным неучтеным фактором здесь является сопротивление воздуха.

Приведенная ниже модификация модели Галилея принадлежит И. Ньютону¹. Известно, что сила лобового сопротивления воздуха F пропорциональна квадрату скорости, т.е. $F = -\beta v^2$. При этом $\beta = 0.5CS\rho$, где ρ — плотность воздуха; S — площадь поперечного сечения; C — коэффициент лобового сопротивления (для многих задач баллистики $C \approx 0.15$).

Обозначим через F_x и F_y горизонтальную и вертикальную проекции вектора лобового сопротивления. Заметим, что $F_x/F = u/v$, $F_y/F = w/v$, $v = \sqrt{u^2 + w^2}$ (рис. 1.2). Следовательно, $F_x = -\beta u \sqrt{u^2 + w^2}$, $F_y = -\beta w \sqrt{u^2 + w^2}$.

Пусть m — масса снаряда. Тогда в силу второго закона Ньютона справедливы уравнения

$$m \frac{du}{dt} = -\beta u \sqrt{u^2 + w^2}, \quad \frac{dx}{dt} = u, \quad (1.4)$$

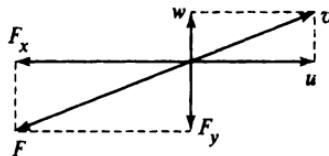


Рис. 1.2

¹ Исаак Ньютон (1643—1727) — английский физик, механик, астроном и математик, заложивший основы современного естествознания.

$$m \frac{dw}{dt} = -mg - \beta w \sqrt{u^2 + w^2}, \quad \frac{dy}{dt} = w, \quad (1.5)$$

которые необходимо дополнить начальными условиями

$$u(0) = v_0 \cos \alpha, \quad w(0) = v_0 \sin \alpha, \quad x(0) = 0, \quad y(0) = 0. \quad (1.6)$$

Полученная модель является более сложной, чем рассмотренная ранее модель (1.1), (1.2), однако она содержит ее как частный случай. Действительно, в случае $\beta = 0$ (сопротивление воздуха отсутствует) уравнения (1.4)–(1.6) и (1.1), (1.2) эквивалентны.

Естественно, что модель (1.4)–(1.6) непригодна для решения задач современной баллистики и реально используемые модели значительно сложнее. ▲

Заметим, что работа по созданию математической модели, как правило, проводится объединенными усилиями специалистов, хорошо знающих предметную область, и математиков, владеющих соответствующими разделами прикладной математики и способных оценить возможность решения возникающих вычислительных задач.

§ 1.2. Основные этапы решения прикладной задачи с применением компьютера

Решение серьезной прикладной задачи с использованием компьютера — довольно длительный и сложный процесс. С определенной степенью условности его можно разбить на ряд последовательных этапов. Выделим следующие этапы

- 1) постановка проблемы;
- 2) выбор или построение математической модели;
- 3) постановка вычислительной задачи;
- 4) предварительный (предмашинный) анализ свойств вычислительной задачи;
- 5) выбор или построение численного метода;
- 6) алгоритмизация и программирование;
- 7) отладка программы;
- 8) счет по программе;
- 9) обработка и интерпретация результатов;
- 10) использование результатов и коррекция математической модели.

1. Постановка проблемы. Первоначально прикладная задача бывает сформулирована в самом общем виде: исследовать некоторое явление, спроектировать устройство, обладающее заданными свойствами, дать прогноз поведения некоторого объекта в определенных условиях и т.д. На данной стадии происходит конкретизация постановки задачи, и пер-

востепенное внимание при этом уделяется выяснению цели исследования. От исследователя требуется глубокое понимание существа задачи и умение сформулировать ее так, чтобы найденное решение было полезным и в то же время могло быть получено с помощью существующих методов и в реальные сроки. Неудачная постановка проблемы может привести к тому, что длительный и дорогостоящий процесс решения задачи завершится получением бесполезных или тривиальных результатов (в этом случае возможно и отсутствие каких-либо результатов).

Этот очень важный и ответственный этап завершается конкретной формулировкой проблемы на языке, принятом в данной предметной области. Знание возможностей, которые дает применение компьютера, может оказать существенное влияние на окончательную формулировку проблемы.

2. Выбор или построение математической модели. Для последующего анализа исследуемого явления или объекта необходимо дать его формализованное описание на языке математики, т.е. построить математическую модель (см. § 1.1). Часто имеется возможность выбора модели среди известных и принятых для описания соответствующих процессов, но нередко требуется и существенная модификация известной модели, а иногда возникает необходимость в построении принципиально новой модели.

Рассматриваемый этап — едва ли не самый важный и трудный. Часто удачный выбор математической модели является решающим шагом к достижению цели. Одна из существенных трудностей такого выбора состоит в объективном противоречии между желанием сделать описание явления как можно более полным (что приводит к усложнению модели) и необходимостью иметь достаточно простую модель (чтобы была возможность реализовать ее на компьютере). Важно, чтобы сложность математической модели соответствовала сложности поставленной проблемы. Если поставленных целей можно достигнуть, используя более простую математическую модель, то ей и следует отдать предпочтение. Как правило, полезно иметь несколько упрощенных вариантов принимаемой модели. Заметим, что грамотное упрощение модели — непростая задача, однако анализ упрощенных моделей весьма полезен в течение всего процесса решения задачи. Такие упрощенные модели часто позволяют ответить на многие принципиальные вопросы и понять основные закономерности поведения более сложной модели.

3. Постановка вычислительной задачи. На основе принятой математической модели формулируют вычислительную задачу (или ряд таких задач). Анализируя результаты ее решения, исследователь предполагает получить ответы на интересующие его вопросы.

4. Предварительный анализ свойств вычислительной задачи. На этом этапе проводят предварительное (предмашинное) исследование свойств вычислительной задачи. Большое внимание уделяют анализу корректности ее постановки, т.е. выяснению вопросов существования и единственности решения, а также исследованию устойчивости решения задачи к погрешностям входных данных (эти вопросы более подробно рассматриваются в гл. 3). Такое исследование, как правило, относится к компетенции профессиональных математиков. Тем не менее, решающему задачу полезно быть в курсе современного состояния названных проблем, уметь самостоятельно проводить простейшие исследования.

К сожалению, для многих имеющих практическую ценность задач их строгое исследование в полной постановке провести не удается и к решению приступают без детального анализа математических свойств этих задач. Это нежелательная, но вынужденная мера, так как в прикладных исследованиях существенное значение имеют конкретные (часто — весьма сжатые) сроки получения результата. На этом этапе полезным оказывается изучение упрощенных постановок задачи. Иногда для них удается провести исследование, позволяющее понять основные особенности исходной вычислительной задачи. Особую ценность имеют различные аналитические решения; они оказываются полезными не только для анализа явления, но и как основа для тестовых испытаний на этапе отладки программы.

5. Выбор или построение численного метода. Для решения вычислительной задачи на компьютере требуется использование численных методов.

Часто решение прикладной задачи сводится к последовательному решению стандартных вычислительных задач, для которых разработаны эффективные численные методы. В этой ситуации происходит либо выбор среди известных методов, либо их адаптация к особенностям решаемой задачи. Однако если возникающая вычислительная задача является новой, то не исключено, что для ее решения не существует готовых методов. Построение численного метода для такой задачи может оказаться очень трудной проблемой и потребовать привлечения специалиста по вычислительной математике. Умение различать отмеченные две ситуации необходимо, и наличие его уже говорит об определенной квалификации в области вычислительных методов.

Для решения одной и той же вычислительной задачи обычно может быть использовано несколько методов. Необходимо знать особенности этих методов, критерии, по которым оценивается их качество, чтобы выбрать метод, позволяющий решить проблему наиболее эффективным образом. Здесь выбор далеко не однозначен. Он существенно зависит

от требований, предъявляемых к решению, от имеющихся в наличии ресурсов, от доступной для использования вычислительной техники и т.д.

Возникающим на этом этапе вопросам и посвящена большая часть данной книги.

6. Алгоритмизация и программирование. Как правило, выбранный на предыдущем этапе численный метод содержит только принципиальную схему решения задачи, не включающую многие детали, без которых невозможна реализация метода на компьютере. Необходима подробная детализация всех этапов вычислений, для того чтобы получить реализуемый на компьютере алгоритм. Составление программы сводится к переводу этого алгоритма на выбранный язык программирования. Заметим, что в настоящее время для вычислительных задач наиболее широко используются алгоритмические языки СИ++ и ФОРТРАН.

В книге значительное место уделяется алгоритмам (в гл. 3 обсуждаются их общие свойства и критерии оценки качества) и практически не рассматриваются вопросы собственно программирования. Конечно, алгоритмизация и программирование очень тесно связаны. Более того, практика показывает, что небольшие, на первый взгляд, различия в программах могут привести к значительным различиям в их эффективности. Тем не менее, вопрос разработки качественного программного продукта мы не затрагиваем (этому предмету посвящено большое число пособий). Подчеркнем лишь, что большинство пользователей справедливо предполагает строить свои программы из готовых модулей и использовать стандартные программы, реализующие те или иные алгоритмы. Разумеется, отсутствие в библиотеке стандартных программ той или иной программы не должно быть непреодолимым препятствием.

7. Отладка программы. На этом этапе с помощью компьютера выявляют и исправляют ошибки в программе.

Как правило, начинающий пользователь компьютера убежден, что ошибок в составленной им программе нет или же они могут быть легко обнаружены и исправлены. Однако совершенно неожиданно для него отладка программы и доведение ее до рабочего состояния нередко оказывается длительным и весьма трудоемким процессом. Приобретая определенный опыт в составлении и отладке сравнительно сложных программ, пользователь убеждается в справедливости популярного афоризма: «В любой программе есть по крайней мере одна ошибка».

Таким образом, наличие в программах ошибок — вполне нормальное и закономерное явление, поэтому подготовку к отладке следует начинать уже на этапе алгоритмизации и программирования. Заметим, что эффективность отладки самым существенным образом зависит от общей методики разработки программ.

После устранения ошибок программирования необходимо провести тщательное тестирование программы — проверку правильности ее работы на специально отобранных тестовых задачах, имеющих известные решения.

8. Счет по программе. На этом этапе происходит решение задачи на компьютере по составленной программе в автоматическом режиме. Этот процесс, в ходе которого входные данные с помощью компьютера преобразуются в результат, называют *вычислительным процессом*. Как правило, счет повторяется многократно с различными входными данными для получения достаточно полной картины зависимости от них решения задачи.

Первые полученные результаты тщательно анализируются, для того чтобы убедиться в правильности работы программы и пригодности выбранного метода решения. Счет по программе продолжается несколько секунд, минут или часов. Именно быстротечность этого этапа порождает распространенную иллюзию о возможности решать важные прикладные задачи на компьютере в очень короткое время. В действительности же, конечно, необходимо принимать во внимание весь цикл от постановки проблемы до использования результатов. Для серьезных задач часто полезные результаты получаются только в результате многолетней работы.

9. Обработка и интерпретация результатов. Полученные в результате расчетов на компьютере выходные данные, как правило, представляют собой большие массивы чисел. Начинающий пользователь часто пытается вывести эти массивы на печать, чтобы «потом провести их анализ». Обычно первый же опыт анализа распечатки, содержащий сотни тысяч чисел, приводит к пониманию того, что эта работа непосильна для человека и следует постараться возложить ее на компьютер.

Часто первоочередной интерес представляет лишь небольшая часть полученной информации (например, значения одной из функций в выделенных точках) или даже некоторая грубая интегральная характеристика (максимальное или минимальное значение, оценка энергии системы и т.д.).

Для того чтобы исследователь мог воспользоваться результатами расчетов, их необходимо представить в виде компактных таблиц, графиков или в иной удобной для восприятия форме. При этом следует, максимально использовать возможности компьютера для подготовки такой информации и ее представления с помощью печатающих и графических выходных устройств.

Для правильной интерпретации результатов расчетов и оценки их достоверности от исследователя требуется глубокое знание существа

решаемой инженерной задачи, ясное представление об используемой математической модели и понимание (хотя бы в общих чертах) особенностей применяемого вычислительного метода.

Вопросы обработки и интерпретации результатов вычислений будут затронуты при рассмотрении конкретных вычислительных методов и алгоритмов.

10. Использование результатов и коррекция математической модели. Завершающий этап состоит в использовании результатов расчетов в практической деятельности, иначе говоря, во внедрении результатов. Не стоит огорчаться, если большинство полученных сначала результатов окажется бесполезным. Действительно полезные для практики результаты являются плодом серьезной целенаправленной работы, в процессе которой цикл решения задачи повторяется неоднократно.

Очень часто анализ результатов, проведенной на этапе их обработки и интерпретации, указывает на несовершенство используемой математической модели и необходимость ее коррекции. В таком случае математическую модель модифицируют (при этом она, как правило, усложняется) и начинают новый цикл решения задачи.

§ 1.3. Вычислительный эксперимент

Создание математических моделей и решение прикладных задач с применением компьютера требует выполнения большого объема работ (см. § 1.1, 1.2). Нетрудно заметить аналогию с соответствующими работами, проводимыми при организации натурных экспериментов: составление программы экспериментов, создание экспериментальной установки, выполнение контрольных экспериментов, проведение серийных опытов, обработка экспериментальных данных и их интерпретация и т.д. Однако вычислительный эксперимент проводится не над реальным объектом, а над его математической моделью, и роль экспериментальной установки играет компьютер, оснащенный специально разработанным пакетом прикладных программ. В связи с этим естественно рассматривать проведение больших комплексных расчетов при решении прикладных и научных задач как *вычислительный эксперимент*, а описанную в предыдущем параграфе последовательность этапов решения как один его цикл.

Широкое применение компьютеров в математическом моделировании, разработанная теория и значительные практические результаты позволяют говорить о вычислительном эксперименте как о новой технологии и методологии научных и прикладных исследований. Серьезное внедрение вычислительного эксперимента в прикладные исследования

лишь начинается, но там, где оно происходит реально (в авиационной и космической промышленности), его плоды весьма весомы.

Отметим некоторые достоинства вычислительного эксперимента по сравнению с натурным. Вычислительный эксперимент, как правило, дешевле физического. В этот эксперимент можно легко и безопасно вмешиваться. Его можно повторить еще раз (если в этом есть необходимость) и прервать в любой момент. В ходе этого эксперимента можно смоделировать условия, которые нельзя создать в лаборатории.

Заметим, что в ряде случаев проведение натурного эксперимента затруднено (а иногда и невозможно), так как изучаются быстропротекающие процессы, исследуются труднодоступные или вообще пока недоступные объекты. Часто проведение полномасштабного натурного эксперимента сопряжено с губительными или непредсказуемыми последствиями (ядерная война, поворот сибирских рек) или с опасностью для жизни или здоровья людей. Нередко требуется исследование и прогнозирование результатов катастрофических явлений (авария ядерного реактора АЭС, глобальное потепление климата, землетрясение). В этих случаях вычислительный эксперимент может стать основным средством исследования. Заметим, что с его помощью оказывается возможным прогнозировать свойства новых, еще не созданных конструкций и материалов на стадии их проектирования.

Существенным недостатком вычислительного эксперимента является то, что применимость его результатов ограничена рамками принятой математической модели.

Конечно, вычислительный эксперимент никогда не сможет полностью заменить натурный, и будущее за их разумным сочетанием. Действительно, построение математической модели основано на результатах наблюдений, опыта, а достоверность ее выводов проверяется с помощью критерия — практики.

Для прикладных задач характерно наличие значительного числа параметров (конструктивных, технологических и др.). Создание нового изделия или технологического процесса предполагает выбор среди большого числа альтернативных вариантов, а также оптимизацию по ряду параметров. Поэтому в ходе вычислительного эксперимента расчеты проводятся многократно с разными значениями входных параметров. Для получения нужных результатов с требуемой точностью и в приемлемые сроки необходимо, чтобы на расчет каждого варианта тратилось минимальное время. Именно поэтому при создании программного обеспечения так важно использовать эффективные численные методы.

Разработка программного обеспечения вычислительного эксперимента в конкретной области прикладной деятельности приводит к созданию крупного программного комплекса. Он состоит из связанных между

собой прикладных программ и системных средств, включающих средства, предоставляемые пользователю для управления ходом вычислительного эксперимента, обработки и представления его результатов. Такой комплекс программ иногда называют *проблемно-ориентированным пакетом прикладных программ*.

§ 1.4. Дополнительные замечания

1. Математическое моделирование является основой современной методологии решения прикладных задач, и его роль объективно возрастает в связи с необходимостью решения все более сложных прикладных проблем.

2. Эффективность применения компьютера в той или иной области науки и техники тем выше, чем совершеннее ее математические модели. В то же время использование компьютера для исследования каких-либо процессов часто служит серьезным стимулом для создания новых математических моделей и детального изучения этих процессов другими методами.

3. Само по себе применение компьютера не позволяет решить прикладную задачу, а лишь дает в руки исследователя мощный инструмент познания. Использование компьютера не только не освобождает от необходимости глубоко осмыслить решаемую проблему, но и заставляет уделять постановке задачи гораздо больше внимания.

4. Вычислительный эксперимент не противоречит натурному эксперименту и классическому математическому анализу прикладной задачи, а, напротив, находится с ними в органическом единстве.

5. Для успешного применения метода математического моделирования с использованием компьютера необходимо гармоническое владение всеми его составляющими. В настоящем пособии основное внимание уделено вычислительным методам и алгоритмам, анализу свойств вычислительных задач, интерпретации получаемых результатов.

6. Дополнительную информацию о методологии современного математического моделирования и концепции вычислительного эксперимента можно получить из следующих источников: [58, 71, 79, 96, 99, 11*].

7. В настоящее время большое количество высококачественных программ и пакетов программ распространяется бесплатно. Например, через Internet по адресу <http://www.netlib.org> доступно программное обеспечение, предназначенное для решения широкого спектра прикладных задач. Кроме того, по адресу <http://www.acm.org/pubs/calgo> можно найти 866 тщательно разработанных программ, реализующих эффективные вычислительные алгоритмы.

Глава 2

ВВЕДЕНИЕ В ЭЛЕМЕНТАРНУЮ ТЕОРИЮ ПОГРЕШНОСТЕЙ

§ 2.1. Источники и классификация погрешностей результата численного решения задачи

Для правильного понимания подходов и критериев, используемых при решении прикладной задачи с применением компьютера, очень важно с самого начала признать, что получить точное значение решения практически невозможно и не в этом цель вычислений. Получаемое на компьютере решение u^* почти всегда (за исключением некоторых весьма специальных случаев) содержит погрешность, т.е. является приближенным. Невозможность получения точного решения следует уже из ограниченной разрядности компьютера.

Наличие погрешности решения обусловлено рядом весьма глубоких причин. Перечислим их.

1. Математическая модель является лишь приближенным описанием реального процесса. Характеристики процесса, вычисленные в рамках принятой модели, заведомо отличаются от истинных характеристик, причем их погрешность зависит от степени адекватности модели реальному процессу.

2. Исходные данные, как правило, содержат погрешности, поскольку они либо получаются в результате экспериментов (измерений), либо являются результатом решения некоторых вспомогательных задач.

3. Применяемые для решения задачи методы в большинстве случаев являются приближенными. Найти решение возникающей на практике задачи в виде конечной формулы возможно только в отдельных, очень упрощенных ситуациях.

4. При вводе исходных данных в компьютер, выполнении арифметических операций и выводе результатов на печать производятся округления.

Пусть u — точное значение величины, вычисление которой является целью поставленной задачи. Соответствующая первым двум из указанных причин погрешность $\delta_u u$ называется *неустранимой погрешностью*. Такое название вызвано тем, что принятие математической модели и задание исходных данных вносит в решение погрешность, которая

не может быть устранена далее. Единственный способ уменьшить эту погрешность — перейти к более точной математической модели и задать более точные исходные данные.

Погрешность $\delta_m u$, источником которой является метод решения задачи, называется *погрешностью метода*, а погрешность $\delta_b u$, возникающая из-за округлений при вводе, выводе и вычислениях, — *вычислительной погрешностью*. Таким образом, полная погрешность результата решения задачи на компьютере $\delta u = u - u^*$ складывается из трех составляющих: неустранимой погрешности, погрешности метода и вычислительной погрешности, т.е. $\delta u = \delta_h u + \delta_m u + \delta_b u$.

Будем далее исходить из предположения, что математическая модель фиксирована и входные данные задаются извне, так что повлиять на значение величины $\delta_h u$ в процессе решения задачи действительно нельзя. Однако это совсем не означает, что предварительные оценки значения неустранимой погрешности не нужны. Достоверная информация о порядке величины $\delta_h u$ позволяет осознанно выбрать метод решения задачи и разумно задать его точность. Желательно, чтобы погрешность метода была в 2—10 раз меньше неустранимой погрешности. Большее значение $\delta_m u$ ощутимо снижает точность результата, меньшее — обычно требует увеличения затрат на вычисления, практически уже не влияя на значение полной погрешности. Иногда характер использования результата таков, что вполне допустимо, чтобы погрешность $\delta_m u$ была сравнима с $\delta_h u$ или даже несколько превышала ее.

Значение вычислительной погрешности (при фиксированных модели, входных данных и методе решения) в основном определяется характеристиками используемого компьютера. Желательно, чтобы погрешность $\delta_b u$ была хотя бы на порядок меньше погрешности метода и совсем не желательна ситуация, когда она существенно ее превышает.

Умение анализировать погрешности при решении прикладной задачи и соблюдать между ними разумный компромисс позволяет существенно экономить используемые ресурсы и является признаком высокой квалификации.

§ 2.2. Приближенные числа.

Абсолютная и относительная погрешности

В § 2.1 было отмечено, что числа, получаемые при решении на компьютере прикладных задач, как правило, являются приближенными. Следовательно, вопрос о *точности* результатов, т.е. о мере их уклонения от истинных значений, в теории и практике методов вычислений

приобретает особое значение. Начнем его рассмотрение с введения основных понятий элементарной теории погрешностей.

Условимся относительно обозначений, которые в дальнейшем будут использоваться при сравнении величин. Кроме привычных знаков « $=$ », « \neq », « $<$ », « \leq », будем использовать знаки приближенного равенства « \approx » и приближенного неравенства « \lesssim ». Когда положительные величины a и b являются величинами одного порядка (т.е. $10^{-1} \lesssim a/b \lesssim 10$), будем использовать обозначение $a \sim b$. Если же a много меньше b , то будем писать $a \ll b$, что эквивалентно соотношению $a/b \ll 1$.

1. Абсолютная и относительная погрешности. Пусть a — точное (вообще говоря, неизвестное) значение некоторой величины, a^* — известное приближенное значение той же величины (*приближенное число*). *Ошибкой или погрешностью* приближенного числа a^* называют разность $a - a^*$ между точным и приближенным значениями.

Простейшей количественной мерой погрешности является *абсолютная погрешность*

$$\Delta(a^*) = |a - a^*|. \quad (2.1)$$

Однако по значению абсолютной погрешности далеко не всегда можно сделать правильное заключение о качестве приближения. Действительно, если $\Delta(a^*) = 0.1$, то следует ли считать погрешность большой или нужно признать ее малой? Ответ существенным образом зависит от принятых единиц измерения и масштабов величин. Если $a \approx 0.3$ то скорее всего точность приближения невелика; если же $a \approx 3 \cdot 10^8$, то следует признать точность очень высокой. Таким образом, естественно соотнести погрешность величины и ее значение, для чего вводится понятие *относительной погрешности* (при $a \neq 0$)

$$\delta(a^*) = \frac{|a - a^*|}{|a|} = \frac{\Delta(a^*)}{|a|}. \quad (2.2)$$

Использование относительных погрешностей удобно, в частности, тем, что они не зависят от масштабов величин и единиц измерения. Заметим, что для приведенного выше примера $\delta(a^*) \approx 0.33 = 33\%$ в первом случае и $\delta(a^*) \approx 0.33 \cdot 10^{-9} = 0.33 \cdot 10^{-7}\%$ во втором.

Так как значение a неизвестно, то непосредственное вычисление величин $\Delta(a^*)$ и $\delta(a^*)$ по формулам (2.1), (2.2) невозможно. Более реальная и часто поддающаяся решению задача состоит в получении оценок погрешности вида

$$|a - a^*| \leq \bar{\Delta}(a^*), \quad (2.3)$$

$$\frac{|a - a^*|}{|a|} \leq \bar{\delta}(a^*), \quad (2.4)$$

где $\bar{\Delta}(a^*)$ и $\bar{\delta}(a^*)$ — величины, которые мы будем называть *верхними границами* (или просто *границами*) *абсолютной* и *относительной погрешностей*.

Если величина $\bar{\Delta}(a^*)$ известна то неравенство (2.4) будет выполнено, если положить

$$\bar{\delta}(a^*) = \frac{\bar{\Delta}(a^*)}{|a|}. \quad (2.5)$$

Точно так же если величина $\bar{\delta}(a^*)$ известна, то следует положить

$$\bar{\Delta}(a^*) = |a|\bar{\delta}(a^*). \quad (2.6)$$

Поскольку значение a неизвестно, при практическом применении формулы (2.5), (2.6) заменяют приближенными равенствами

$$\bar{\delta}(a^*) \approx \frac{\bar{\Delta}(a^*)}{|a^*|}, \quad \bar{\Delta}(a^*) \approx |a^*|\bar{\delta}(a^*). \quad (2.7)$$

Замечание. В литературе по методам вычислений широко используется термин «точность». Принято говорить о точности входных данных и решения, о повышении и снижении точности вычислений и т.д. Мы также будем использовать эту терминологию, за которой скрывается довольно простой смысл. Точность в качественных рассуждениях обычно выступает как противоположность погрешности, хотя для количественного их измерения используются одни и те же характеристики (например, абсолютная и относительная погрешности). Точное значение величины — это значение, не содержащее погрешности. Повышение точности воспринимается как уменьшение погрешности, а снижение точности — как увеличение погрешности. Часто используемая фраза «требуется найти решение с заданной точностью ε » означает, что ставится задача о нахождении приближенного решения, принятая мера погрешности которого не превышает заданного значения ε . Вообще говоря, следовало бы говорить об абсолютной точности и относительной точности, но часто этого не делают, считая, что из контекста ясно, как измеряется погрешность.

2. Правила записи приближенных чисел. Пусть приближенное число a^* задано в виде конечной десятичной дроби:

$$a^* = \alpha_n\alpha_{n-1}\dots\alpha_0 \cdot \beta_1\beta_2\dots\beta_m.$$

Значащими цифрами числа a^* называют все цифры в его записи, начиная с первой ненулевой слева.

Пример 2.1. У чисел $a^* = 0.0\cancel{1}03$ и $a^* = 0.0\cancel{1}03000$ значащие цифры подчеркнуты. Первое число имеет три, а второе — шесть значащих цифр. ▲

Значащую цифру числа a^* называют *верной*, если абсолютная погрешность числа не превосходит единицы разряда, соответствующего этой цифре.

Пример 2.2. Если $\bar{\Delta}(a^*) = 2 \cdot 10^{-6}$, то число $a^* = 0.0\cancel{1}03000$ имеет четыре верные значащие цифры (они подчеркнуты). ▲

Следует отметить, что широко распространенной ошибкой при записи приближенных чисел является отбрасывание последних значащих нулей (даже если они представляют собой верные цифры).

Замечание. Верная цифра приближенного числа, вообще говоря, не обязана совпадать с соответствующей цифрой в записи точного числа. Таким образом, термин «верная цифра» не следует понимать буквально (см. пример 2.3).

Пример 2.3. Пусть $a = 1.00000$, $a^* = 0.\cancel{9}99999$. Тогда $\Delta(a^*) = 0.00001$ и у числа a^* все подчеркнутые цифры — верные, хотя они и не совпадают с соответствующими цифрами числа a . ▲

Количество верных значащих цифр числа тесно связано со значением его относительной погрешности. Приведенные ниже утверждения позволяют в дальнейшем связывать точность числа с количеством его верных значащих цифр и трактовать потерю точности как потерю верных цифр.

Предложение 2.1. 1°. *Если число a^* содержит N верных значащих цифр, то справедливо неравенство*

$$\delta(a^*) \leq (10^{N-1} - 1)^{-1} \approx 10^{-N+1}.$$

2°. *Для того чтобы число a^* содержало N верных значащих цифр, достаточно чтобы было выполнено неравенство*

$$\delta(a^*) \leq (10^N + 1)^{-1} \approx 10^{-N}.$$

3°. *Если число a^* имеет ровно N верных значащих цифр, то*

$$10^{-N-1} \leq \delta(a^*) \leq 10^{-N+1}$$

и, таким образом, $\delta(a^) \sim 10^{-N}$.* ■

Пример 2.4. Что можно сказать об относительной погрешности числа a^* , если известно, что оно содержит три верные значащие цифры?

В силу утверждения 1° имеем $\delta(a^*) \leq 10^{-2} = 1\%$. ▲

Пример 2.5. С какой относительной точностью следует найти число a^* , чтобы верными оказались шесть его значащих цифр?

Из утверждения 2° следует, что достаточно найти a^* с относительной точностью $\varepsilon \approx 10^{-6}$. ▲

Заметим, что границы абсолютной и относительной погрешностей принято записывать с одной или двумя значащими цифрами. Большая точность в записи этих величин, как правило, не имеет смысла, так как обычно они являются довольно грубыми оценками истинных значений погрешностей, и кроме того, для практического использования часто бывает достаточно знать только их порядок.

Пример 2.6. Информация о погрешности вида $\delta(a^*) \approx 0.288754 \cdot 10^{-5}$ практически равнозначна информации $\delta(a^*) \approx 3 \cdot 10^{-6}$, причем последняя вызывает больше доверия. Скорее всего, вполне удовлетворительной в данном случае является и запись $\delta(a^*) \sim 10^{-6}$. ▲

Неравенство (2.3) эквивалентно двойному неравенству

$$a^* - \bar{\Delta}(a^*) \leq a \leq a^* + \bar{\Delta}(a^*)$$

и поэтому тот факт, что число a^* является приближенным значением числа a с верхней границей абсолютной погрешности $\bar{\Delta}(a^*)$ [с абсолютной точностью $\varepsilon = \bar{\Delta}(a^*)$], принято записывать в виде

$$a = a^* \pm \bar{\Delta}(a^*).$$

Как правило, числа a^* и $\bar{\Delta}(a^*)$ указывают с одинаковым числом цифр после десятичной точки.

Пример 2.7. Пусть для числа a известны приближенное значение $a^* = 1.648$ и граница абсолютной погрешности $\bar{\Delta}(a^*) = 0.002832$. Тогда можно записать $a = 1.648 \pm 0.003$. Записи вида $a = 1.648 \pm 0.002832$ или $a = 1.648 \pm 0.1$ являются неестественными. ▲

Из неравенства (2.4) следует, что значение a заключено примерно между $a^*(1 - \bar{\delta}(a^*))$ и $a^*(1 + \bar{\delta}(a^*))$. Поэтому тот факт, что число a^* является приближенным значением числа a с границей относительной погрешности $\bar{\delta}(a^*)$ (с относительной точностью $\varepsilon = \bar{\delta}(a^*)$), принято записывать в виде $a = a^*(1 \pm \bar{\delta}(a^*))$.

Пример 2.8. Оценим точность часто используемого в простейших расчетах приближения $\pi^* = 3\frac{14}{7}$ к числу π . Поскольку $\pi = 3.14159\dots$, то $\pi - \pi^* = 0.00159\dots$ Следовательно, можно принять $\bar{\Delta}(\pi^*) = 0.0016$ и $\bar{\delta}(\pi^*) \approx 0.0016/3\frac{14}{7} \approx 0.00051 = 0.051\%$. Итак, $\pi = 3.14(1 \pm 0.051\%)$. ▲

Замечание. Если число a^* приводится в качестве результата без указания значения погрешности, то принято считать, что все его значащие цифры являются верными. Начинающий пользователь часто слишком доверяет выводимым из компьютера цифрам, предполагая, что вычислительная машина придерживается того же соглашения. Однако это совсем не так: число может быть выведено с таким количеством значащих цифр, сколько потребует программист заданием соответствующего формата. Как правило, среди этих цифр только небольшое число первых окажутся верными, а возможно верных цифр нет совсем. Анализировать результаты вычислений и определять степень их достоверности совсем непросто. Одна из целей изучения вычислительных методов и состоит в достижении понимания того, что можно и чего нельзя ожидать от результатов, полученных на компьютере.

3. Округление. Часто возникает необходимость в *округлении* числа a , т.е. в замене его другим числом a^* с меньшим числом значащих цифр. Возникающая при такой замене погрешность называется *погрешностью округления*.

Существует несколько способов округления числа до n значащих цифр. Наиболее простой из них — *усечение* состоит в отбрасывании всех цифр, расположенных справа от n -й значащей цифры. Более предпочтительным является *округление по дополнению*. В простейшем варианте это правило округления состоит в следующем. Если первая слева из отбрасываемых цифр меньше 5, то сохраняемые цифры остаются без изменения. Если же она больше либо равна 5, то в младший сохраняемый разряд добавляется единица.

Абсолютное значение погрешности округления при округлении по дополнению не превышает половины единицы разряда, соответствующего последней оставляемой цифре, а при округлении усечением — единицы того же разряда.

Пример 2.9. При округлении числа $a = 1.72631$ усечением до трех значащих цифр получится число $a^* = 1.72$, а при округлении по дополнению — число $a^* = 1.73$. ▲

Границы абсолютной и относительной погрешностей принято всегда округлять в сторону увеличения.

Пример 2.10. Округление значений $\bar{\Delta}(a^*) = 0.003721$ и $\bar{\delta}(a^*) = 0.0005427$ до двух значащих цифр дает значения $\bar{\Delta}(a^*) = 0.0038$ и $\bar{\delta}(a^*) = 0.00055$. ▲

§ 2.3. Погрешность арифметических операций над приближенными числами

Исследуем влияние погрешностей исходных данных на погрешность результатов арифметических операций. Пусть a^* и b^* — приближенные значения чисел a и b . Какова соответствующая им величина неустранимой погрешности результата?

Предложение 2.2. Абсолютная погрешность алгебраической суммы (суммы или разности) не превосходит суммы абсолютных погрешностей слагаемых, т.е.

$$\Delta(a^* \pm b^*) \leq \Delta(a^*) + \Delta(b^*). \quad (2.8)$$

Доказательство. Имеем

$$\begin{aligned} \Delta(a^* \pm b^*) &= |(a \pm b) - (a^* \pm b^*)| = |(a - a^*) \pm (b - b^*)| \leq \\ &\leq \Delta(a^*) + \Delta(b^*). \blacksquare \end{aligned}$$

Следствие. В силу неравенства (2.8) естественно положить

$$\bar{\Delta}(a^* \pm b^*) = \bar{\Delta}(a^*) + \bar{\Delta}(b^*). \quad (2.9)$$

Оценим относительную погрешность алгебраической суммы.

Предложение 2.3. Пусть a и b — ненулевые числа одного знака. Тогда справедливы неравенства

$$\delta(a^* + b^*) \leq \delta_{\max}, \quad \delta(a^* - b^*) \leq v \delta_{\max}, \quad (2.10)$$

где $\delta_{\max} = \max\{\delta(a^*), \delta(b^*)\}$, $v = |a + b|/|a - b|$.

Доказательство. Используя формулу (2.2) и неравенство (2.8), имеем

$$\begin{aligned} |a \pm b| \delta(a^* \pm b^*) &= \Delta(a^* \pm b^*) \leq \Delta(a^*) + \Delta(b^*) = \\ &= |a| \delta(a^*) + |b| \delta(b^*) \leq (|a| + |b|) \delta_{\max} = |a + b| \delta_{\max}. \end{aligned}$$

Из полученного неравенства сразу следуют оценки (2.10). ■

Следствие. В силу неравенств (2.10) естественно положить

$$\bar{\delta}(a^* + b^*) = \bar{\delta}_{\max}, \quad \bar{\delta}(a^* - b^*) = v \bar{\delta}_{\max}, \quad (2.11)$$

где $\bar{\delta}_{\max} = \max\{\bar{\delta}(a^*), \bar{\delta}(b^*)\}$, $v = |a + b|/|a - b|$.

Первое из равенств (2.11) означает, что при суммировании чисел одного знака не происходит потери точности, если оценивать точность в относительных единицах. Совсем иначе обстоит дело при вычитании

чисел одного знака. Здесь граница относительной погрешности возрастает в $v > 1$ раз и возможна существенная потеря точности. Если числа a и b близки настолько, что $|a + b| \gg |a - b|$, то $v \gg 1$ и не исключена полная или почти полная потеря точности. Когда это происходит, говорят о том, что произошла *катастрофическая потеря точности*.

Пример 2.11. Пусть решается прикладная задача, в которой окончательный результат y вычисляется по формуле $y = 1 - x$ с помощью предварительно определяемого значения x . Предположим, что найденное приближение $x^* = 0.999997$ к значению x содержит шесть верных значащих цифр. Тогда $y^* = 1 - 0.999997 = 0.000003$ и в процессе вычисления оказались потерянными пять верных цифр. Если же учесть, что $\delta(x^*) \sim 0.0001\%$, а $\delta(y^*) \sim 33\%$, то следует признать, что произошла катастрофическая потеря точности.

Подчеркнем, что здесь виновником «катастрофы» является не операция вычитания, а предложенный метод решения задачи, где окончательный результат получается с помощью вычитания двух близких чисел. Выполнение этой операции лишь делает очевидным то, что действительно полезная информация о значении y уже оказалась потерянной до вычитания. Если нет другого варианта расчета, то для получения приемлемого результата следовало бы предварительно вычислить x с существенно большим числом верных знаков, учитывая, что пять старших значащих цифр при вычитании будут потеряны. ▲

Итак, получаем следующий важный вывод. При построении численного метода решения задачи следует избегать вычитания близких чисел одного знака. Если же такое вычитание неизбежно, то следует вычислять аргументы с повышенной точностью, учитывая ее потерю примерно в $v = |a + b|/|a - b|$ раз.

Предложение 2.4. Для относительных погрешностей произведения и частного приближенных чисел верны оценки

$$\delta(a^* b^*) \leq \delta(a^*) + \delta(b^*) + \delta(a^*) \delta(b^*), \quad (2.12)$$

$$\delta(a^*/b^*) \leq \frac{\delta(a^*) + \delta(b^*)}{1 - \delta(b^*)}, \quad (2.13)$$

в последней из которых считается, что $\delta(b^*) < 1$.

Доказательство. Выполним следующие преобразования:

$$\begin{aligned} |ab| \delta(a^* b^*) &= \Delta(a^* b^*) = |ab - a^* b^*| = \\ &= |(a - a^*)b + (b - b^*)a - (a - a^*)(b - b^*)| \leq \\ &\leq |b| \Delta(a^*) + |a| \Delta(b^*) + \Delta(a^*) \Delta(b^*) = |ab| (\delta(a^*) + \delta(b^*) + \delta(a^*) \delta(b^*)). \end{aligned}$$

Разделив обе части полученного неравенства на $|ab|$, выведем оценку (2.12).

Для вывода второй оценки предварительно заметим, что

$$|b^*| = |b + (b^* - b)| \geq |b| - \Delta(b^*) = |b|(1 - \delta(b^*)).$$

Тогда

$$\begin{aligned} \delta(a^*/b^*) &= \frac{|a/b - a^*/b^*|}{|a/b|} = \frac{|ab^* - ba^*|}{|ab^*|} = \\ &= \frac{|a(b^* - b) + b(a - a^*)|}{|ab^*|} \leq \frac{|a|\Delta(b^*) + |b|\Delta(a^*)}{|ab^*(1 - \delta(b^*))|} = \frac{\delta(b^*) + \delta(a^*)}{1 - \delta(b^*)}. \blacksquare \end{aligned}$$

Следствие. Если $\bar{\delta}(a^*) \ll 1$ и $\bar{\delta}(b^*) \ll 1$, то для оценки границ относительных погрешностей можно использовать следующие приближенные равенства:

$$\bar{\delta}(a^*b^*) \approx \bar{\delta}(a^*) + \bar{\delta}(b^*), \quad \bar{\delta}(a^*/b^*) \approx \bar{\delta}(a^*) + \bar{\delta}(b^*). \quad (2.14)$$

Именно равенства (2.14) чаще всего и используют для практической оценки погрешности.

Итак, выполнение арифметических операций над приближенными числами, как правило, сопровождается потерей точности. Единственная операция, при которой потеря не происходит, — это сложение чисел одного знака. Наибольшая потеря точности может произойти при вычитании близких чисел одного знака.

§ 2.4. Погрешность функции

1. Погрешность функции многих переменных. Пусть $f(x) = f(x_1, x_2, \dots, x_m)$ — дифференцируемая в области G функция m переменных, вычисление которой производится при приближенно заданных значениях аргументов $x_1^*, x_2^*, \dots, x_m^*$. Такая ситуация возникает, например всякий раз, когда на компьютере производится расчет по формуле. Важно знать, каково значение неустранимой погрешности, вызванной тем, что вместо значения $y = f(x)$ в действительности вычисляется значение $y^* = f(x^*)$, где $x^* = (x_1^*, x_2^*, \dots, x_m^*)$.

Введем обозначения: пусть $[x, x^*]$ — отрезок¹, соединяющий точки x и x^* , и $f'_{x_j} = \partial f / \partial x_j$.

¹ Отрезком, соединяющим точки x и x^* в m -мерном пространстве, называется множество точек вида $\alpha x + (1 - \alpha)x^*$, $0 \leq \alpha \leq 1$.

Предложение 2.5. Для абсолютной погрешности значения $y^* = f(x^*)$ справедлива следующая оценка:

$$\Delta(y^*) \leq \sum_{j=1}^m \max_{\tilde{x} \in [x, x^*]} |f'_{x_j}(\tilde{x})| \Delta(x_j^*). \quad (2.15)$$

Доказательство. Оценка (2.15) вытекает из формулы конечных приращений Лагранжа¹:

$$f(x) - f(x^*) = \sum_{j=1}^m f'_{x_j}(\tilde{x})(x_j - x_j^*), \quad \tilde{x} \in [x, x^*]. \blacksquare$$

Следствие. Если $x^* \approx x$, то в силу оценки (2.15) можно положить

$$\bar{\Delta}(y^*) \approx \sum_{j=1}^m |f'_{x_j}(x^*)| \bar{\Delta}(x_j^*), \quad (2.16)$$

$$\bar{\Delta}(y^*) \approx \sum_{j=1}^m |f'_{x_j}(x)| \bar{\Delta}(x_j^*). \quad (2.17)$$

Равенство (2.16) удобно для практических оценок, а равенством (2.17) мы воспользуемся в дальнейшем для теоретических построений.

Из (2.16), (2.17) вытекают приближенные равенства для оценки границ относительных погрешностей:

$$\bar{\delta}(y^*) \approx \sum_{j=1}^m v_j^* \bar{\delta}(x_j^*), \quad \bar{\delta}(y^*) \approx \sum_{j=1}^m v_j \bar{\delta}(x_j^*). \quad (2.18)$$

$$v_j^* = \frac{|x_j^*| |f'_{x_j}(x^*)|}{|f(x^*)|}, \quad v_j = \frac{|x_j| |f'_{x_j}(x)|}{|f(x)|}. \quad (2.19)$$

Пример 2.12. Пусть корни квадратного уравнения $x^2 + bx + c = 0$ вычисляются при значениях коэффициентов $b \approx 10^3$, $c \approx 1$. Каково влияние погрешностей задания коэффициентов на точность вычисляемых значений?

Воспользуемся явными формулами для корней:

$$x_1 = f(b, c) = (-b - \sqrt{d})/2, \quad x_2 = g(b, c) = (-b + \sqrt{d})/2,$$

где $d = b^2 - 4c$.

¹ Жозеф Луи Лагранж (1736—1813) — французский математик, механик и астроном. Один из создателей математического анализа, вариационного исчисления, классической аналитической механики.

Заметим, что $x_1 \cdot x_2 = c$. Тогда при заданных значениях коэффициентов получим:

$$\sqrt{d} \approx \sqrt{10^6 - 4} \approx 10^3, \quad x_1 \approx (-10^3 - 10^3)/2 = -10^3,$$

$$x_2 = c/x_1 \approx 1/(-10^3) = -10^{-3}.$$

Далее, имеем:

$$f'_b = (-1 - b/\sqrt{d})/2 \approx (-1 - 10^3/10^3)/2 = -1, \quad f'_c = 1/\sqrt{d} \approx 10^{-3},$$

$$g'_b = (-1 + b/\sqrt{d})/2 = (b - \sqrt{d})/(2\sqrt{d}) = -x_2/\sqrt{d} \approx 10^{-6},$$

$$g'_c = -1/\sqrt{d} \approx -10^{-3}.$$

Применяя первую из формул (2.19), для корня x_1 находим:

$$v_{1,1}^* = |b^*| |f'_b| / |x_1^*| \approx 1, \quad v_{1,2}^* = |c^*| |f'_c| / |x_1^*| \approx 10^{-6}.$$

Аналогично для корня x_2 имеем:

$$v_{2,1}^* = |b^*| |g'_b| / |x_2^*| \approx 1, \quad v_{2,2}^* = |c^*| |g'_c| / |x_2^*| \approx 1.$$

Таким образом,

$$\bar{\delta}(x_1^*) \approx \bar{\delta}(b^*) + 10^{-6}\bar{\delta}(c^*), \quad \bar{\delta}(x_2^*) \approx \bar{\delta}(b^*) + \bar{\delta}(c^*). \quad (2.20)$$

Следовательно, точность первого корня практически определяется только точностью задания коэффициента b , в то время как коэффициент c может быть задан очень грубо. Для второго корня влияние погрешностей в задании коэффициентов b и c практически одинаково. ▲

2. Погрешность функции одной переменной. Формулы для границ погрешностей функции $f(x)$ одной переменной являются частным случаем формул (2.16)–(2.18) при $m = 1$:

$$\bar{\Delta}(y^*) \approx |f'(x^*)| \bar{\Delta}(x^*), \quad \bar{\Delta}(y^*) \approx |f'(x)| \bar{\Delta}(x^*); \quad (2.21)$$

$$\bar{\delta}(y^*) \approx v^* \bar{\delta}(x^*), \quad \bar{\delta}(y^*) \approx v \bar{\delta}(x^*), \quad (2.22)$$

$$\text{где } v^* = \frac{|x^*| |f'(x^*)|}{|f(x^*)|}, \quad v = \frac{|x| |f'(x)|}{|f(x)|}.$$

3. Погрешность неявной функции. Нередко приходится сталкиваться с ситуацией, когда функция $y = f(x)$, где $x = (x_1, x_2, \dots, x_m)$ задается не явным образом, а как решение нелинейного уравнения $F(y, x) = 0$, т.е. неявно. Если для такой неявной функции воспользоваться известными формулами вычисления частных производных

$$f'_{x_j}(x) = - \left. \frac{F'_{x_j}(y, x)}{F'_y(y, x)} \right|_{y=f(x)}, \quad j = 1, 2, \dots, m, \quad (2.23)$$

то исследование неустранимой погрешности неявной функции сразу же сводится к рассмотренному выше случаю.

Пример 2.13. Для проведенного в примере 2.12 исследования совсем не обязательно было выписывать явные формулы для корней. В этом случае величины $x_1 = f(b, c)$ и $x_2 = g(b, c)$ можно рассматривать как неявные функции, заданные уравнением $F(x, b, c) = 0$, где $F = x^2 + bx + c$.

Здесь $F'_x = 2x + b$, $F'_b = x$, $F'_c = 1$. Следовательно,

$$v_{1,i}^* = \frac{|b| |F'_b/F'_x|}{|x|} \Bigg|_{x=x_i^*, b=10^3, c=1} \approx \frac{10^3}{|2x_i^* + 10^3|}$$

$$v_{2,i}^* = \frac{|c| |F'_c/F'_x|}{|x|} \Bigg|_{x=x_i^*, b=10^3, c=1} \approx \frac{1}{|(2x_i^* + 10^3)x_i^*|}.$$

Вычисления при $i = 1$, $x_1^* \approx -10^3$ и $i = 2$, $x_2^* \approx -10^{-3}$ дают те же значения коэффициентов $v_{1,1}^*$, $v_{1,2}^*$, $v_{2,1}^*$, $v_{2,2}^*$, что и в примере 2.12, а следовательно, те же формулы (2.20). ▲

§ 2.5. Особенности машинной арифметики

Знание основных особенностей машинной арифметики необходимо для грамотного использования компьютеров при решении научно-технических задач. Пользователь, не учитывающий эти особенности, вряд ли может рассчитывать на высокую точность и эффективность вычислений. Невнимание к ним часто приводит к неверным результатам. Подчеркнем, что в основе причин появления вычислительной погрешности лежит сам способ представления чисел в компьютерах.

1. Системы счисления. Принятый способ записи чисел состоит в представлении их упорядоченным набором цифр. В привычной нам десятичной позиционной системе счисления вещественное число x представляют последовательностью символов, которая начинается со знака (+ или -) и продолжается цепочкой десятичных цифр α_i и β_j , разделенных десятичной точкой.

$$x = \pm \alpha_n \dots \alpha_1 \alpha_0 . \beta_1 \beta_2 \dots \beta_m \dots \quad (2.24)$$

Здесь каждой позиции (разряду), которую занимает цифра относительно десятичной точки, отвечает определенная степень числа 10.

По существу, равенство (2.24) представляет собой принятую форму сокращенной записи разложения числа x в сумму по степеням 10:

$$x = \pm(\alpha_n \cdot 10^n + \dots + \alpha_1 \cdot 10^1 + \alpha_0 \cdot 10^0 + \beta_1 \cdot 10^{-1} + \\ + \beta_2 \cdot 10^{-2} + \dots + \beta_m \cdot 10^{-m} + \dots).$$

Пример 2.14. Запись $x = 20.5$ означает, что $x = 2 \cdot 10^1 + 0 \cdot 10^0 + 5 \cdot 10^{-1}$. ▲

Для представления чисел в компьютере также используют позиционные системы счисления, однако основаниями систем служат, как правило, степени числа 2. Это вызвано способом хранения чисел в устройствах памяти компьютера, каждое из которых можно рассматривать как набор однотипных элементов, способных находиться только в одном из двух возможных устойчивых состояний — «включен» или «выключен». Эти состояния интерпретируются соответственно как 0 или 1 — значение двоичной цифры. Наиболее распространены системы счисления с основанием 2 (базисная двоичная система счисления), 8 и 16.

Игнорируя некоторые малосущественные детали, будем считать, что все вычислительные машины работают в двоичной системе счисления. В ней вещественное число x по-прежнему записывается в виде (2.24), однако α_i и β_j — уже двоичные цифры (0 или 1). В этом случае число x разлагается в сумму по степеням числа 2:

$$x = \pm(\alpha_n \cdot 2^n + \dots + \alpha_1 \cdot 2^1 + \alpha_0 \cdot 2^0 + \beta_1 \cdot 2^{-1} + \beta_2 \cdot 2^{-2} + \dots + \beta_m \cdot 2^{-m} + \dots).$$

Пример 2.15. Запишем число $x = 20.5$ в двоичной системе счисления. Для этого разложим его в сумму по степеням числа 2:

$$x = 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 0 \cdot 2^0 + 1 \cdot 2^{-1}.$$

Опуская степени числа 2, получаем $x = (10100.1)_2$. Здесь нижний индекс 2 указывает на основание системы счисления. ▲

Для хранения числа в памяти компьютера отводится поле стандартной длины (*машинное слово*), в котором число записывают в виде последовательности двоичных цифр. По форме представления, способу хранения и реализации арифметических операций существенно различаются два типа используемых в компьютерах чисел: целые числа и вещественные числа.

2. Представление целых чисел.

Целое число n представляют в виде

$$n = \pm(\alpha_L \cdot 2^L + \dots + \alpha_1 \cdot 2^1 + \alpha_0 \cdot 2^0), \quad (2.25)$$

где L — некоторое стандартное для компьютера целое число, α_i — двоичные цифры. Всего для хранения числа n отводят $s = L + 2$ разрядов (один из них для хранения знака).

Из представления (2.25) видно, что максимальное по модулю целое число, представимое в компьютере, есть $n_{\max} = 2^L + \dots + 2^1 + 2^0 = 2^{L+1} - 1$. Обычно оно не очень велико. Например, при стандартном формате записи целых чисел в компьютерах $s = 32$ и $n_{\max} = 2^{31} - 1 \approx 2 \cdot 10^9$.

Операции сложения, вычитания и умножения над целыми числами реализованы так, что если результат не превышает по модулю число n_{\max} , то он получается точным. Отметим, однако, следующую неприятную особенность. Если модуль результата превышает n_{\max} , то на большинстве компьютеров эта ситуация не доводится до сведения пользователя, происходит присвоение результату некоторого значения (меньшего n_{\max} по модулю) и вычисления продолжаются далее.

3. Представление вещественных чисел. В современных компьютерах для вещественных чисел принята форма представления с плавающей точкой, когда каждое число представляют в виде

$$x = \pm(\gamma_1 \cdot 2^{-1} + \gamma_2 \cdot 2^{-2} + \dots + \gamma_t \cdot 2^{-t})2^p. \quad (2.26)$$

Здесь $\gamma_1, \gamma_2, \dots, \gamma_t$ — двоичные цифры.

Как правило, число x нормализуется так, чтобы $\gamma_1 = 1$, и поэтому в памяти компьютера хранятся только значащие цифры соответствующего нормализованного числа. Число $\mu = \pm(\gamma_1 \cdot 2^{-1} + \gamma_2 \cdot 2^{-2} + \dots + \gamma_t \cdot 2^{-t})$ называется *мантиссой* числа x . Количество t цифр, которое отводится для записи мантиссы, называемое *разрядностью мантиссы*, зависит от конструктивных особенностей конкретной вычислительной машины, но всегда является конечным. В представлении (2.26) p — целое число, называемое *двоичным порядком*. Порядок также записывают как двоичное целое число $p = \pm(\sigma_1\sigma_{l-1}\dots\sigma_0)_2$, для хранения которого в машинном слове отводится $l+2$ двоичных разрядов. На рис. 2.1 схематически представлена структура машинного слова для хранения вещественного нормализованного числа.

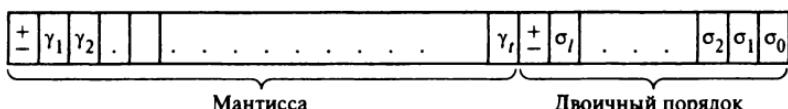


Рис. 2.1

Пример 2.16. Представим число $x = 20.5$ в виде двоичного нормализованного числа. Так как $x = (10100.1)_2$ (см. пример 2.15), то, перемещая двоичную точку на пять позиций влево, получаем $x = (0.101001)_2 \cdot 2^5$. ▲

Длительное время различные компьютеры отличались друг от друга числом разрядов, отводимых для хранения мантиссы и порядка, способом округления и имели различные правила машинной арифметики. Однако в настоящее время большинство компьютеров (в том числе — все персональные компьютеры) конструируются в соответствии с разработанным в 1985 году IEEE — стандартом двоичной арифметики (IEEE Floating Point Standard). При стандартном способе записи нормализованного числа x в IEEE арифметике под хранение мантиссы отводится 24 разряда (включая знак), а под хранение двоичного порядка — 8 разрядов. Поскольку для нормализованного числа всегда $\gamma_1 = 1$, то необходимости в хранении первого значащего разряда нет, и сэкономленный разряд используется для хранения еще одного двоичного разряда мантиссы.

Поскольку $\gamma_1 = 1$, то для мантиссы нормализованного числа справедливы оценки $0.5 \leq |\mu| < 1$. В то же время для представления порядка используется конечное число $l + 1$ двоичных цифр и поэтому $-(2^{l+1} - 2) = p_{\min} \leq p \leq p_{\max} = 2^{l+1} - 1$. Таким образом, для представимых на компьютере нормализованных чисел имеем $0 < X_0 \leq |x| < X_\infty$, где

$X_0 = 2^{p_{\min}-1}$, $X_\infty = 2^{p_{\max}}$. Числа X_0 и X_∞ иногда называют *порогом машинного нуля* и *порогом переполнения*.

На компьютерах, удовлетворяющих стандарту IEEE, в арифметике обычной точности $p_{\min} = -126$, $p_{\max} = 127$ и поэтому $X_0 = 2^{-127} \approx 10^{-38}$, $X_\infty = 2^{127} \approx 10^{38}$.

Кроме нормализованных чисел IEEE арифметика включает в себя еще и *субнормальные (денормализованные) числа*. Субнормальные числа — очень малые числа из диапазона $(-X_0, X_0)$, которые возникают в процессе вычислений и не могут быть представлены в нормализованном виде. Они также представляются в виде (2.26), но с $\gamma_1 = 0$ и минимально возможным показателем $p = p_{\min}$. Напомним, что в IEEE арифметике обычной точности $X_0 \approx 10^{-38}$.

Пример 2.17. Число $x = 2^{-130}$ не может быть представлено в IEEE арифметике обычной точности как нормализованное число. Оно представляется как субнормальное число вида $x = (0.00010 \dots 0)_2 \cdot 2^{-126}$. ▲

Заметим, что нуль — особое ненормализуемое число. Оно представляется в IEEE арифметике как +0 и -0. Числа +0 и -0 считаются равными, но в некоторых случаях они могут проявлять себя по разному. В представлении (2.26) нуль имеет нулевую мантиссу и показатель степени $p = p_{\min} - 1$, т.е. $\pm 0 = \pm(0.000 \dots 0)_2 \cdot 2^{p_{\min} - 1}$.

На основании имеющихся сведений о представлении чисел в компьютере можно сделать ряд важных выводов.

1°. В компьютере представимы не все числа, а лишь конечный набор рациональных чисел специального вида. Эти числа образуют *представимое множество* компьютера. Для всех остальных чисел x возможно лишь их приближенное представление с ошибкой, которую принято называть *ошибкой представления* (или *ошибкой округления*). Обычно приближенное представление числа x в компьютере обозначают как $x^* = \text{fl}(x)$ ¹. В IEEE арифметике округление производится по дополнению, поэтому для нормализованных чисел граница относительной погрешности представления равна единице первого отброшенного разряда мантиссы, т.е. $\bar{\delta}(x^*) = \varepsilon_m = 2^{-t}$ (порядок числа не влияет на относительную погрешность представления). Величина ε_m играет в вычислениях на компьютере фундаментальную роль; ее называют *относительной точностью* компьютера, а также *машинной точностью* (или *машинным эпсилоном*). Всюду в дальнейшем ε_m — это относительная точность компьютера. Заметим, что значение этой величины определяется разрядностью мантиссы и способом округления.

Важно с самого начала иметь четкое представление о том, что почти наверняка в представимом множестве компьютера нет числа y , являющегося решением поставленной задачи. Лучшее, что можно попытаться сделать, — это найти его представление $y^* = \text{fl}(y)$ с относительной точностью порядка ε_m .

Полезно отметить, что среди представимых в компьютере чисел нет не только ни одного иррационального (в том числе и таких важных постоянных, как π , e , $\sqrt{2}$), но и даже такого широко используемого в вычислениях числа, как 0.1. Дело в том, что двоичная запись числа 0.1 является бесконечной периодической дробью: $0.1 =$

¹ fl — начальные буквы английского слова floating — «плавающий».

$= (0.0001100110011\dots)_2$. Поэтому это число всегда представляется в компьютере приближенно, с погрешностью, вызванной необходимостью округления.

2°. Диапазон изменения чисел в компьютере ограничен. Нормализованные числа расположены в интервалах $(-X_\infty, -X_0)$ и (X_0, X_∞) , где значения X_0 и X_∞ определяются числом разрядов, отведенных для хранения двоичного порядка.

3°. Все числа x , по модулю большие X_∞ , рассматриваются как *машинная бесконечность* и представляются в IEEE арифметике символами $-\infty$ или $+\infty$. Получение числа x такого, что $|x| > X_0$, называют *переполнением*. Получение числа x такого, что $|x| < X_0$, называют *исчезновением порядка* или *антитереполнением*. Обычно при антипереполнении генерируется субнормальное число или полагается $\text{fl}(x) = 0$ и вычисления продолжаются.

Замечание. Не следует смешивать машинную точность ϵ_m , порог машинного нуля X_0 и минимальное представимое на компьютере положительное субнормальное число x_{\min} . Это совершенно различные числа, причем $x_{\min} \ll X_0 \ll \epsilon_m$. В IEEE арифметике обычной точности $x_{\min} \approx 10^{-43}$, $X \approx 10^{-38}$ и $\epsilon_m \approx 10^{-7}$.

4°. На машинной числовой оси (рис. 2.2) числа расположены неравномерно. Плотность их возрастает по мере приближения к нулю и падает с удалением от нуля. Чтобы убедиться в этом, заметим, что расстояние от одного нормализованного числа $x = \mu \cdot 2^p$ до другого ближайшего нормализованного равно единице последнего разряда мантиссы, умноженной на 2^p , т.е. равно 2^{p-t} . Так как t фиксировано, то расстояние уменьшается с уменьшением значения порядка p и возрастает с увеличением p . Субнормальные числа распределены на интервале $(-X_0, X_0)$ равномерно (см. рис. 2.2), расстояние между соседними субнормальными числами равно $2^{p_{\min}-t}$



Рис. 2.2

4. Арифметические операции над числами с плавающей точкой. Правила выполнения арифметических операций в двоичной системе счисления чрезвычайно просты и легко реализуются на компьютере. Однако в силу ограниченной разрядности мантиссы операции сложения, вычитания, умножения и деления над представимыми в компьютере вещественными числами не могут быть реализованы точно. Дело в том, что арифметические операции над числами, мантиссы которых содержат t разрядов, приводят, как правило, к результатам, содержащим более t разрядов. Округление результата до t разрядов и служит главным источником погрешности. Для того чтобы отличать машинные арифметические операции от идеальных математических операций $+$, $-$, \times , $:$, будем обозначать их через \oplus , \ominus , \otimes , \circ . Игнорируя несущественные детали, можно считать, что результат машинной арифметической операции совпадает с результатом точного выполнения той же операции с погрешностью, приближенно равной погрешности округления. Таким образом,

$$\bar{\Delta}(a \oplus b) \approx |a + b| \varepsilon_m, \quad \bar{\Delta}(a \ominus b) \approx |a - b| \varepsilon_m,$$

$$\bar{\Delta}(a \otimes b) \approx |a \times b| \varepsilon_m, \quad \bar{\Delta}(a \circ b) \approx |a : b| \varepsilon_m.$$

Конечно, в некоторых ситуациях округление может отсутствовать. Например, полезно знать, что умножение и деление числа на целую степень числа 2 выполняется на компьютере точно, так как в этом случае мантисса не меняется.

Пример 2.18. Рассмотрим гипотетический компьютер, в котором числа представляются всего лишь с шестью двоичными разрядами мантиссы, а округление производится по дополнению. Пусть на таком компьютере вычисляются сумма и произведение двух представимых на ней чисел $a = 20.5 = (10100.1)_2$ и $b = 1.75 = (1.11)_2$. Производим вычисления в двоичной арифметике:

$$a + b = (10100.1)_2 + (1.11)_2 = (10110.01)_2,$$

$$a \times b = (10100.1)_2 \times (1.11)_2 = (100011.111)_2.$$

После округления до шести значащих цифр получим $a \oplus b = (10110.1)_2 = 22.5$, $a \otimes b = (100100.)_2 = 36$. Очевидно, что эти результаты отличаются от точных значений $a + b = 22.25$, $a \times b = 35.875$. ▲

Заметим, что машинные арифметические операции обладают иными свойствами, нежели обычные математические операции. Например, не выполняется известное правило арифметики «от перемены мест слагаемых сумма не меняется». Покажем это на примере.

Пример 2.19. Пусть вычисления производятся на компьютере из примера 2.17, причем $a = (1.)_2$, $b = c = (0.000001)_2$. Тогда $a + b = (1.000001)_2$ и после округления имеем $a \otimes b = (1.00001)_2$. Далее, $(a \otimes b) + c = (1.000011)_2$ и после округления получим $(a \otimes b) \otimes c = (1.00010)_2$. Сложение в ином порядке дает $c \otimes b = (0.000010)_2$, $(c \otimes b) \otimes a = (1.00001)_2$. Таким образом, $(a \otimes b) \otimes c \neq (c \otimes b) \otimes a$. \blacktriangle

5. Двойная точность. На компьютерах, удовлетворяющих IEEE стандарту, при вычислениях с обычной точностью, нормализованные числа расположены в диапазоне $10^{-38} \leq |x| \leq 10^{38}$, вполне достаточном для большинства приложений. Однако разрядность мантиссы t невелика и поэтому $\epsilon \sim 10^{-7}$; в десятичной арифметике это эквивалентно тому, что мантисса содержит семь десятичных цифр.

На большинстве компьютеров возможна реализация арифметических действий над числами, разрядность мантиссы которых примерно вдвое превосходит стандартную разрядность t . Это приводит к существенному повышению машинной точности. Например на компьютерах, удовлетворяющих стандарту IEEE, для представления чисел двойной точности отводится 53 разряда для хранения мантиссы (включая знак) и 11 разрядов для хранения порядка. Здесь $X_0 \approx 10^{-308}$, $X_\infty \approx 10^{308}$ и $\epsilon_m \sim 10^{-16}$. Для сравнения напомним, что обычная точность этих компьютеров есть $\epsilon_m \sim 10^{-7}$, так что следует говорить не об удвоении точности, а о повышении точности на много порядков.

В IEEE арифметике есть еще режим расширенной двойной точности, в котором для представления мантиссы отводится 65 разрядов (включая знак), а для представления порядка — 15 разрядов. В этом режиме $X_0 \approx 10^{-4964}$, $X_\infty \approx 10^{4964}$ и $\epsilon_m \sim 10^{-19}$. Подчеркнем, что использование режимов двойной или расширенной двойной точности отнюдь не позволяет ликвидировать погрешность округления, а только лишь уменьшает ее значение.

6. Вычисление машинного ϵ . Для приближенного вычисления величины ϵ_m удобно воспользоваться следующим определением. Машинное эпсилон — это минимальное из представимых на компьютере положительных чисел ϵ , для которых $1 + \epsilon > 1$.

Величину ϵ_m можно оценить непосредственно в ходе вычислительного процесса. Для этого достаточно включить в программу фрагмент, реализующий следующий метод. Полагая $\epsilon^{(0)} = 1$, следует вычислять последовательно $\epsilon^{(1)} = 0.5\epsilon^{(0)}$, $\epsilon^{(2)} = 0.5\epsilon^{(1)}$, ..., $\epsilon^{(n)} = 0.5\epsilon^{(n-1)}$, ...,

запоминая получающиеся значения $\varepsilon^{(n)}$ в памяти компьютера и проверяя каждый раз выполнение неравенства $1 + \varepsilon^{(n)} > 1$. Как только при некотором n окажется, что $1 + \varepsilon^{(n)} = 1$, следует положить $\varepsilon_m = \varepsilon^{(n-1)}$ и перейти к следующему этапу вычислений. Хотя полученное таким способом значение может отличаться от ε_m в два раза, обычно оно используется так, что эта погрешность не имеет значения.

Пример 2.20. Покажем, что $\varepsilon_m = 2^{-6} = (0.000001)_2$ — машинное эпсилон для компьютера из примера 2.18. В самом деле, $1 + \varepsilon_m = (1.000001)_2$ и после округления имеем $1 + \varepsilon_m = (1.00001)_2 > 1$. Если же к 1 добавить любое положительное $\varepsilon < \varepsilon_m$, то в седьмом разряде результата будет стоять нуль и после округления получим $1 + \varepsilon = 1$. ▲

Всюду в дальнейшем, приводя конкретные числовые примеры, мы отказываемся от использования двоичной арифметики. Десятичная арифметика привычнее, а основные закономерности поведения ошибок округления не зависят от основания используемой системы. В большинстве расчетов, которые будут приведены для иллюстрации поведения ошибок округления, имитируется выполнение вычислений на гипотетическом компьютере, имеющем шесть десятичных разрядов мантиссы и производящем округление по дополнению. Будем называть этот компьютер *б-разрядным десятичным компьютером*. Для него $\varepsilon_m = 5 \cdot 10^{-7}$, так что по точности он сравним со стандартным компьютером (при вычислениях с одинарной точностью).

7. Дополнительные сведения об IEEE арифметике. Напомним, что кроме нормализованных чисел IEEE арифметика включает в себя субнормальные (денормализованные) числа. В отличие от нормализованных чисел в представлении (2.26) для субнормальных чисел всегда $\gamma_1 = 0$ и $p = p_{\min}$. Субнормальные числа расположены в интервале $(-X_0, X_0)$. Расстояние между двумя соседними субнормальными числами постоянно. В арифметике обычной точности оно составляет $2^{-150} \approx 10^{-45}$, а в арифметике двойной точности $2^{-1075} \approx 10^{-324}$.

Одним из приятных следствий наличия субнормальных чисел является тот факт, что в IEEE арифметике равенство $x \Theta y = 0$ возможно тогда и только тогда, когда $x = y$.

Множество нормализованных и субнормальных чисел в IEEE арифметике дополнено символами $+\infty$, $-\infty$ и *NaN* (Not a Number — НеЧисло). Символы $\pm\infty$ генерируются в случае, если при выполнении очередной операции возникает переполнение. Символы $\pm\infty$ представля-

ются в виде (2.26) с нулевой мантиссой и показателем $p = p_{\max} + 1$, т.е. в виде

$$\pm\infty = \pm(0.000 \dots 0) \cdot 2^{p_{\max} + 1}.$$

Символы $\pm\infty$ ведут себя в соответствии со следующими естественными правилами: $+\infty + \infty = +\infty$, $(-1) \times (+\infty) = -\infty$, $\frac{x}{\pm\infty} = 0$ для всякого числа x с плавающей точкой, $\frac{\pm\infty}{x} = \pm\infty$ для всякого положительного числа x с плавающей точкой, $+\infty \oplus x = +\infty$ и т.д.

Символ Nan генерируется в случае, если результат выполняемой операции не определен. Например $\frac{0}{0} = Nan$, $\sqrt{-1} = Nan$, $0 \times (+\infty) = Nan$, $+\infty + (-\infty) = Nan$, $\frac{\infty}{\infty} = Nan$, $Nan \oplus x = Nan$, $Nan \ominus Nan = Nan$ и т.д.

Nan числа представляются в виде (2.26) с ненулевой мантиссой и с показателем $p = p_{\max} + 1$.

Таким образом, IEEE арифметика замкнута относительно арифметических операций. Каждая выполняемая на компьютере операция имеет результатом либо число с плавающей точкой, либо один из символов $+\infty$, $-\infty$, Nan . Информация о возникновении символов $\pm\infty$ или Nan доводится до сведения пользователя выставлением *индикатора особого случая* (exception flag), но работа программы при этом не прерывается. Это позволяет разрабатывать более надежные программы. Любые два Nan числа считаются различными. Поэтому простейший способ убедиться в том, что $x = Nan$ состоит в проверке свойства $x \neq x$.

Замечание. Казалось бы, результат выполнения операции 0^0 в IEEE арифметике должен получить значение Nan , но в действительности принято соглашение, что $0^0 = 1$. Это соглашение очень полезно при работе с многочленами.

§ 2.6. Дополнительные замечания

- Более подробно приближенные числа и погрешности арифметических операций над ними изложены в известном учебном пособии [35]. Здесь же приведены правила подсчета оставляемых значащих цифр, которые рекомендуется применять при массовых «ручных» вычислениях без точного учета погрешностей.

2. Дополнительную информацию об особенностях арифметических операций над числами с плавающей точкой можно найти, например в [23, 54, 102, 106].

3. В последнее время в практике вычислений в качестве меры погрешности приближенного числа a^* часто используют величину $\frac{|a - a^*|}{1 + |a|}$, объединяющую в себе черты абсолютной и относительной погрешностей. Она близка к $\Delta(a^*)$ при $|a| \ll 1$ и практически совпадает с $\delta(a^*)$ при $|a| \gg 1$.

4. Информацию о IEEE стандарте можно найти в [54, 122, 5*, 15*, 20*, 24*].

5. Для уменьшения погрешностей округления иногда используется *рациональная арифметика*. В ней каждое рациональное число представляется парой целых чисел (с увеличенным по сравнению со стандартным числом разрядов) — числителем и знаменателем. Арифметические операции над такими числами производятся как над обычными рациональными дробями.

Рациональная арифметика реализована во многих пакетах символьных вычислений, например в пакетах Mathematica и MACSYMA.

6. К настоящему времени разработан *интервальный анализ*, дающий возможность заставить компьютер следить за совершамыми ошибками округления, чтобы затем получить гарантированные оценки их накопления.

Глава 3

ВЫЧИСЛИТЕЛЬНЫЕ ЗАДАЧИ, МЕТОДЫ И АЛГОРИТМЫ. ОСНОВНЫЕ ПОНЯТИЯ

§ 3.1. Корректность вычислительной задачи

1. Постановка вычислительной задачи. Под *вычислительной задачей* будем понимать одну из трех задач, которые возникают при анализе математических моделей: прямую задачу, обратную задачу или задачу идентификации (см. § 1.2). Слово «вычислительная» подчеркивает, что основные усилия будут направлены на то, чтобы найти (вычислить) ее решение.

Будем считать, что постановка задачи включает в себя задание *множества допустимых входных данных X и множества возможных решений Y* . Цель вычислительной задачи состоит в нахождении решения $y \in Y$ по заданному входному данному $x \in X$. Для простоты понимания достаточно ограничиться рассмотрением задач, в которых входные данные и решение могут быть только числами, наборами чисел (векторами, матрицами, последовательностями) и функциями. Предположим, что для оценки величин погрешностей приближенных входных данных x^* и приближенного решения y^* введены абсолютные и относительные погрешности $\Delta(x^*)$, $\Delta(y^*)$, $\delta(x^*)$, $\delta(y^*)$, а также их границы $\bar{\Delta}(x^*)$, $\bar{\Delta}(y^*)$, $\bar{\delta}(x^*)$, $\bar{\delta}(y^*)$. Определения этих величин когда x и y — числа, были даны в § 2.2. В тех случаях, когда входные данные или решение не являются числами, эти характеристики погрешностей также можно ввести естественным образом; мы будем это делать по мере необходимости.

2. Определение корректности задачи. Анализ важнейших требований, предъявляемых к различным прикладным задачам, приводит к понятию корректности математической задачи, которое было впервые сформулировано Ж. Адамаром¹ и развито затем И.Г. Петровским². Вычислительная задача называется *корректной* (по Адамару—Петровскому), если выполнены следующие три требования:

- 1) ее решение $y \in Y$ существует при любых входных данных $x \in X$;

¹ Жак Адамар (1865—1963) — французский математик.

² Иван Георгиевич Петровский (1901—1973) — российский математик.

- 2) это решение единственno;
 3) решение устойчиво по отношению к малым возмущениям входных данных.

В том случае, когда хотя бы одно из этих требований не выполнено, задача называется *некорректной*.

Существование решения вычислительной задачи — естественное требование к ней. Отсутствие решения может свидетельствовать, например, о непригодности принятой математической модели либо о неправильной постановке задачи. Иногда отсутствие решения является следствием неправильного выбора множества допустимых входных данных X или множества возможных решений Y .

Пример 3.1. Рассмотрим задачу о решении квадратного уравнения

$$x^2 + bx + c = 0. \quad (3.1)$$

Старший коэффициент a считается равным единице; этого всегда можно добиться делением уравнения на a . Если считать входным данным пару коэффициентов b, c и искать решение в множестве вещественных чисел, то существование решений

$$x_1 = (-b - \sqrt{b^2 - 4c})/2, \quad x_2 = (-b + \sqrt{b^2 - 4c})/2 \quad (3.2)$$

будет гарантировано только в том случае, если ограничить множество входных данных коэффициентами, удовлетворяющими условию $b^2 - 4c \geq 0$. Если же расширить множество возможных решений и считать, что корни (3.2) могут принимать комплексные значения, то задача будет иметь решение при любых b, c . ▲

Так как математическая модель не является абсолютно точным отражением реальной ситуации, то даже в случае, когда исходная проблема заведомо имеет решение, соответствующая вычислительная задача может и не оказаться разрешимой. Конечно, такая ситуация говорит о серьезном дефекте в постановке задачи. Иногда отсутствие решения математической задачи приводит к пониманию того, что первоначально сформулированная проблема неразрешима и нуждается в серьезной корректировке.

3. Единственность. Для некоторых вычислительных задач единственность является естественным свойством; для других же решение может и не быть единственным. Например, квадратное уравнение (3.1) имеет два корня (3.2). Как правило, если задача имеет реальное содержание, то неединственность может быть ликвидирована введением дополнительных ограничений на решение (т.е. сужением множества Y). В некоторых случаях проблема снимается тем, что признается целесообразным найти набор всех решений, отвечающих входным данным x ,

и тогда за решение у принимается этот набор. Например, для уравнения (3.1) решением можно назвать пару (x_1, x_2) .

Неединственность решения вычислительной задачи — весьма неприятное свойство. Оно может быть проявлением неправильной постановки исходной прикладной проблемы, неоднозначности ее решения или сигналом о неудачном выборе математической модели.

4. Устойчивость решения. Решение y вычислительной задачи называется *устойчивым по входным данным* x , если оно зависит от входных данных непрерывным образом. Это означает, что для любого $\varepsilon > 0$ существует $\delta = \delta(\varepsilon) > 0$ такое, что всякому исходному данному x^* , удовлетворяющему условию $\Delta(x^*) < \delta$, отвечает приближенное решение y^* , для которого $\Delta(y^*) < \varepsilon$. Таким образом, для устойчивой вычислительной задачи ее решение теоретически можно найти со сколь угодно высокой точностью ε , если обеспечена достаточно высокая точность δ входных данных. Схематическая ситуация изображена на рис. 3.1. Множества тех x^* и y^* , для которых $\Delta(x^*) < \delta$ и $\Delta(y^*) < \varepsilon$ изображены как окрестности точек x и y имеющие радиусы δ и ε . Требование увеличить точность решения приводит автоматически к повышению требований к точности данных; соответствующие окрестности на рис. 3.1 отмечены штрихами.

Неустойчивость решения y означает, что существует такое $\varepsilon_0 > 0$, что какое бы малое $\delta > 0$ ни было задано, найдутся такие исходные данные x^* , что $\Delta(x^*) < \delta$, но $\Delta(y^*) \geq \varepsilon_0$.

Приведем простейшие примеры устойчивых и неустойчивых задач.

Пример 3.2. Задача вычисления корней квадратного уравнения (3.1) устойчива, так как корни (3.2) являются непрерывными функциями коэффициентов b и c . ▲

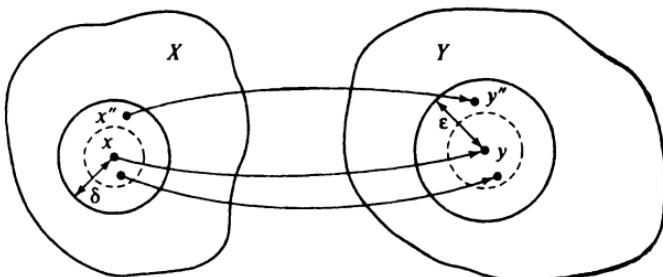


Рис. 3.1

Пример 3.3. Задача о вычислении ранга матрицы в общем случае неустойчива. В самом деле, для матрицы $A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ ранг равен 1, поскольку $\det A = 0$ и существует ненулевой элемент $a_{11} = 1$. Однако сколь угодно малое возмущение коэффициента a_{22} на величину $\varepsilon \neq 0$ приводит к матрице $A_\varepsilon = \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \end{bmatrix}$, для которой $\det A_\varepsilon = \varepsilon \neq 0$ и, следовательно, ранг равен 2. ▲

Пример 3.4. Покажем, что задача вычисления определенного интеграла $I = \int_a^b f(x) dx$ устойчива.

Пусть $f^*(x)$ — приближенно заданная интегрируемая функция и $I^* = \int_a^b f^*(x) dx$. Определим абсолютную погрешность функции f^* с помощью равенства $\Delta(f^*) = \sup_{x \in [a,b]} |f(x) - f^*(x)|$, в котором знак \sup можно заменить на \max , если f и f^* непрерывны. Так как

$$\Delta(I^*) = |I - I^*| = \left| \int_a^b (f(x) - f^*(x)) dx \right| \leq (b - a) \Delta(f^*), \quad (3.3)$$

то для любого $\varepsilon > 0$ неравенство $\Delta(I^*) < \varepsilon$ будет выполнено, если потребовать выполнение условия $\Delta(f^*) < \delta = \varepsilon / (b - a)$. ▲

Пример 3.5. Покажем, что задача вычисления производной $u(x) = f'(x)$ приближенно заданной функции является неустойчивой.

Пусть $f^*(x)$ — приближенно заданная на отрезке $[a, b]$ непрерывно дифференцируемая функция и $u^*(x) = (f^*)'(x)$. Определим абсолютные погрешности с помощью равенств

$$\Delta(f^*) = \max_{[a,b]} |f(x) - f^*(x)|, \quad \Delta(u^*) = \max_{[a,b]} |u(x) - u^*(x)|.$$

Возьмем, например $f^*(x) = f(x) + \alpha \sin(x/\alpha^2)$, где $0 < \alpha \ll 1$. Тогда $u^*(x) = u(x) + \alpha^{-1} \cos(x/\alpha^2)$ и $\Delta(u^*) = \alpha^{-1}$, в то время как $\Delta(f^*) = \alpha$. Таким образом, сколь угодно малой погрешности задания функции f может отвечать сколь угодно большая погрешность производной f' . ▲

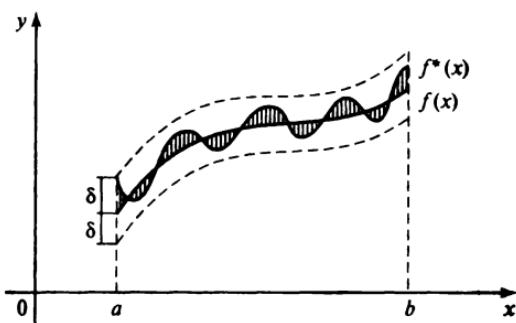


Рис. 3.2

Различие в ситуациях, возникающих при приближенном задании функции f в задачах интегрирования и дифференцирования, отражено на рис. 3.2. Видно, что уменьшение значения $\delta = \Delta(f^*)$ влечет за собой уменьшение погрешности $\Delta(f^*)$ (ее значение не превышает заштрихованной площади), в то время как производные f' и $(f^*)'$ могут отличаться сколь угодно сильно.

Одна и та же задача может оказаться как устойчивой, так и неустойчивой в зависимости от выбора способа вычисления абсолютных погрешностей $\Delta(x^*)$ и $\Delta(y^*)$. В реальных задачах этот выбор определяется тем, в каком смысле должно быть близко приближенное решение к точному и малость какой из мер погрешности входных данных можно гарантировать.

Пример 3.6. Рассмотрим задачу о вычислении суммы сходящегося ряда $S = \sum_{k=0}^{\infty} a_k$ с приближенно заданными слагаемыми $a_k^* \approx a_k$. Если a_k^* определяется таким образом, что гарантируется малость $\Delta(a_k^*)$ для всех k , то для последовательности $a^* = \{a_k^*\}_{k=0}^{\infty}$ естественно положить $\Delta(a^*) = \sup_{k \geq 0} |a_k - a_k^*|$. В такой постановке задача неустойчива. Чтобы убедиться в этом, достаточно положить $a_k^* = a_k + \delta$ (где $\delta > 0$) для $k < N$ и $a_k^* = a_k$ для $k \geq N$. Тогда для суммы ряда $S^* = \sum_{k=0}^{\infty} a_k^*$ имеем $\Delta(S^*) = N\delta$. Следовательно, как бы ни была мала величина $\Delta(a^*) = \delta$, абсолютную погрешность суммы ряда S^* с помощью выбора N можно сделать сколь

угодно большой. Если же положить $a_k^* = a_k + \delta$ для всех k , то сумма ряда S^* вообще станет бесконечной, т.е. ряд станет расходящимся.

В то же время, если можно задавать a_k^* так, чтобы оказалась малой величина $\Delta(a^*) = \sum_{k=0}^{\infty} |a_k - a_k^*|$, то $\Delta(S^*) = \left| \sum_{k=0}^{\infty} (a_k - a_k^*) \right| \leq \Delta(a^*)$ и в такой постановке задача устойчива. ▲

5. Относительная устойчивость решения. Часто требование малости абсолютной погрешности является неоправданным или трудно проверяемым. В таких случаях полезно рассмотреть *относительную устойчивость* решения, определение которой отличается от данного выше определения устойчивости (*абсолютной устойчивости*) только тем, что $\Delta(x^*)$ и $\Delta(y^*)$ заменяются на $\delta(x^*)$ и $\delta(y^*)$ соответственно.

Пример 3.7. Вернемся к задаче вычисления суммы ряда $S = \sum_{k=0}^{\infty} a_k$ из примера 3.6. Предположим, что $a_k \neq 0$ для всех k . Часто можно гарантировать малость величины $\delta(a^*) = \sup_{k \geq 0} \{ |a_k - a_k^*| / |a_k| \}$. Тогда

$\Delta(a_k^*) \leq |a_k| \cdot \delta(a^*)$ и поэтому $\Delta(S^*) \leq \sum_{k=0}^{\infty} |a_k| \cdot \delta(a^*)$. Таким образом,

$$\delta(S^*) \leq \left(\sum_{k=0}^{\infty} |a_k| / \left| \sum_{k=0}^{\infty} a_k \right| \right) \delta(a^*). \quad (3.4)$$

Следовательно, задача вычисления суммы сходящегося ряда S относительно устойчива, если он сходится абсолютно (т.е. сходится ряд $\sum_{k=0}^{\infty} |a_k|$). Если же ряд сходится только условно, т.е. $\sum_{k=0}^{\infty} |a_k| = \infty$, то задача не является относительно устойчивой. ▲

Замечание. Так как для решения вычислительных задач используют компьютеры, точность которых определяется разрядностью мантиссы или эквивалентным значением границы относительной погрешности округления ε_m , то представляется более естественным исследование относительной устойчивости.

6. О некорректных задачах. Длительное время считалось, что некорректные задачи, решения которых неустойчивы, не имеют физического смысла и не представляют ценности для приложений. Однако это мнение оказалось ошибочным. Как выяснилось, многие важные приклад-

ные задачи некорректны. Не вызывает, например, сомнения практическая важность решения некорректных задач дифференцирования и суммирования ряда (см. примеры 3.5., 3.6). К некорректным задачам относятся также обратные задачи геофизики, астрофизики, спектрографии, многие задачи распознавания образов, задачи синтеза и ряд других прикладных задач.

К настоящему времени разработана теория решения многих классов некорректных задач. Важная роль в создании методов решения таких задач принадлежит российским математикам, в первую очередь А.Н. Тихонову¹. Эти методы (*методы регуляризации*) довольно сложны и выходят за рамки данной книги. Мы ограничимся изложением в § 16.7 того, как метод регуляризации может быть применен в случае некорректной задачи вычисления решения интегрального уравнения Фредгольма первого ряда. Для первого знакомства с методами регуляризации можно рекомендовать учебные пособия [37, 97].

§ 3.2. Обусловленность вычислительной задачи

1. Определения. Пусть вычислительная задача корректна (ее решение существует, единственно и устойчиво по входным данным). Теоретически наличие у задачи устойчивости означает, что ее решение может быть найдено со сколь угодно малой погрешностью, если только гарантировать, что погрешности входных данных достаточно малы. Однако на практике погрешности входных данных не могут быть сделаны сколь угодно малыми, точность их ограничена. Даже то что исходные данные нужно будет ввести в компьютер, означает, что их относительная точность будет заведомо ограничена величиной порядка ϵ_m (см. § 2.5). В реальности, конечно, уровень ошибок в исходной информации будет существенно выше. Как же повлияют малые, но конечные погрешности входных данных на решение, как сильно способны они исказить желаемый результат? Для ответа на этот вопрос введем новые понятия.

Под *обусловленностью вычислительной задачи* понимают чувствительность ее решения к малым погрешностям входных данных. Задачу называют *хорошо обусловленной*, если малым погрешностям входных данных отвечают малые погрешности решения, и *плохо обусловленной*, если возможны сильные изменения решения.

Часто оказывается возможным ввести количественную меру степени обусловленности вычислительной задачи — *число обусловленности*. Этую величину можно интерпретировать как коэффициент возможного

¹ Андрей Николаевич Тихонов (1906—1993) — российский математик, один из создателей теории решения некорректных задач.

возрастания погрешностей в решении по отношению к вызвавшим их погрешностям входных данных.

Пусть между абсолютными погрешностями входных данных x и решения y установлено неравенство

$$\Delta(y^*) \leq v_\Delta \Delta(x^*). \quad (3.5)$$

Тогда величина v_Δ называется *абсолютным числом обусловленности*. Если же установлено неравенство

$$\delta(y^*) \leq v_\delta \delta(x^*) \quad (3.6)$$

между относительными погрешностями данных и решения, то величину v_δ называют *относительным числом обусловленности*. В неравенствах (3.5), (3.6) вместо погрешностей Δ и δ могут фигурировать их границы $\bar{\Delta}$ и $\bar{\delta}$. Обычно под числом обусловленности v задачи понимают одну из величин v_Δ или v_δ , причем выбор бывает ясен из смысла задачи. Чаще все же под числом обусловленности понимают относительное число обусловленности. Для *плохо обусловленной задачи* $v \gg 1$. В некотором смысле неустойчивость задачи — это крайнее проявление плохой обусловленности, отвечающее значению $v = \infty$. Конечно, v — это максимальный коэффициент возможного возрастания уровня ошибок, и для конкретных входных данных действительный коэффициент возрастания может оказаться существенно меньше. Однако при выводе оценок (3.5) и (3.6) стремятся к тому, чтобы не завышать значений v_Δ и v_δ и поэтому соотношение $v \gg 1$ все же свидетельствует о реальной возможности существенного роста ошибок. Грубо говоря, если $v \sim 10^N$, где v — относительное число обусловленности, то порядок N показывает число верных цифр, которое может быть утеряно в результате по сравнению с числом верных цифр входных данных.

Каково то значение v , при котором следует признать задачу плохо обусловленной? Ответ на этот вопрос существенно зависит, с одной стороны, от предъявляемых требований к точности решения и, с другой — от уровня обеспечиваемой точности исходных данных. Например, если требуется найти решение с точностью 0.1 %, а входная информация задается с точностью 0.02 %, то уже значение $v = 10$ сигнализирует о плохой обусловленности. Однако (при тех же требованиях к точности результата) гарантия, что исходные данные задаются с точностью не ниже 0.0001 %, означает, что и при $v = 10^3$ задача хорошо обусловлена.

С простым, но очень важным примером плохо обусловленной задачи мы уже познакомились в § 2.3. Это задача вычитания приближенных чисел одного знака. Оценка (2.10) дает для нее значение относительного числа обусловленности $v = |a + b|/|a - b|$. Так как в при-

мере 2.11 имеем $v \approx |1 + x^*| / |1 - x^*| \approx 7 \cdot 10^5$, то потеря пяти верных значащих цифр здесь не представляется удивительной.

2. Примеры плохо обусловленных задач. Первым рассмотрим классический пример, принадлежащий Дж. Уилкинсону¹ [102].

Пример 3.8. Пусть требуется найти корни многочлена

$$P(x) = (x - 1)(x - 2) \dots (x - 20) = x^{20} - 210x^{19} + \dots$$

по заданным значениям его коэффициентов. Из теории известно, что эта задача устойчива. Возьмем коэффициент $\alpha = -210$ при x^{19} и изменим его значение на $\alpha^* = -210 + 2^{-23}$. Как повлияет эта, казалось бы, незначительная погрешность на значения корней? Заметим, что $x_1 = 1$, $x_2 = 2, \dots, x_{20} = 20$ — точные значения корней.

Вычисленные с высокой точностью корни возмущенного многочлена таковы:

$$x_1^* \approx 1.00000, \quad x_2^* \approx 2.00000, \quad x_3^* \approx 3.00000, \quad x_4^* \approx 4.00000,$$

$$x_5^* \approx 5.00000, \quad x_6^* \approx 6.00001, \quad x_7^* \approx 6.99970, \quad x_8^* \approx 8.00727,$$

$$x_9^* \approx 8.91725, \quad x_{10, 11}^* \approx 10.0953 \pm 0.643501i, \quad x_{12, 13}^* \approx 11.7936 \pm 1.65233i,$$

$$x_{14, 15}^* \approx 13.9924 \pm 2.51883i, \quad x_{16, 17}^* \approx 16.7307 \pm 2.81262i,$$

$$x_{18, 19}^* \approx 19.5024 \pm 1.94033i, \quad x_{20}^* \approx 20.8469.$$

Как нетрудно видеть, корни x_1, \dots, x_6 оказались практически нечувствительны к погрешностям в коэффициенте α . В то же время некоторые корни превратились в комплексные и имеют относительные погрешности от 6 до 18 %, несмотря на то, что $\delta(\alpha^*) \approx 6 \cdot 10^{-8} \%$.

В данном случае нетрудно провести анализ чувствительности корней. Пусть $F(x, \alpha) = x^{20} + \alpha x^{19} + \dots$. Будем рассматривать корни x_k как функции параметра α , т.е. $x_k = x_k(\alpha)$. Равенство $F(x_k(\alpha), \alpha) = 0$, выполненное в окрестности точки $\alpha = -210$, задает x_k как неявную функцию от α . Пользуясь второй из формул (2.22) для границы относительной погрешности и применяя формулу (2.23) для производной неявной функции, имеем

$$\bar{\delta}(x_k^*) \approx v_k \bar{\delta}(\alpha^*),$$

¹ Джеймс Х. Уилкинсон (1918—1986) — американский математик, внесший значительный вклад в современное понимание вычислительных методов. О деятельности Джеймса Х. Уилкинсона см. [54].

$$\text{где } v_k = \frac{|\alpha| \cdot |F'_\alpha / F'_x|}{|x|} \Big|_{x=x_k, \alpha=-210}$$

Учитывая, что

$$\frac{\partial F}{\partial \alpha} = x^{19} \quad \text{и} \quad \left. \frac{\partial F}{\partial x} \right|_{x=k, \alpha=-210} = P'(x)|_{x=k} = \sum_{j=1}^{20} \prod_{\substack{i=1 \\ i \neq j}}^{20} (x-i)|_{x=k} = \prod_{\substack{i=1 \\ i \neq k}}^{20} (k-i),$$

получаем

$$v_k = \frac{210k^{18}}{\prod_{\substack{i=1 \\ i \neq k}}^{20} |k-i|}.$$

Вычисление чисел v_k дает следующие значения:

$$\begin{aligned} v_1 &\approx 2 \cdot 10^{-15}, \quad v_2 \approx 9 \cdot 10^{-8}, \quad v_3 \approx 10^{-4}, \quad v_4 \approx 10^{-1}, \quad v_5 \approx 3 \cdot 10^1, \\ v_6 &\approx 2 \cdot 10^3, \quad v_7 \approx 8 \cdot 10^4, \quad v_8 \approx 2 \cdot 10^6, \quad v_9 \approx 2 \cdot 10^7, \quad v_{10} \approx 2 \cdot 10^8, \\ v_{11} &\approx 9 \cdot 10^9, \quad v_{12} \approx 4 \cdot 10^9, \quad v_{13} \approx 10^{10}, \quad v_{14} \approx 2 \cdot 10^{10}, \\ v_{15} &\approx 3 \cdot 10^{10}, \quad v_{16} \approx 3 \cdot 10^{10}, \quad v_{17} \approx 3 \cdot 10^{10}, \quad v_{18} \approx 10^{10}, \\ v_{19} &\approx 3 \cdot 10^9, \quad v_{20} \approx 5 \cdot 10^8, \end{aligned}$$

свидетельствующие о чрезвычайно плохой обусловленности старших корней. ▲

Следует обратить серьезное внимание на то, что задача вычисления корней многочленов высокой степени часто оказывается плохо обусловленной. Поэтому имеет смысл с определенной настороженностью относиться к алгоритмам, составной частью которых является вычисление корней многочленов высокой степени.

К сожалению, эта задача может быть плохо обусловленной и для многочленов невысокой степени, в особенности, если вычисляются кратные корни.

Пример 3.9. Пусть ищется решение уравнения $(x - 1)^4 = 0$ с кратным корнем. Погрешность в младшем коэффициенте, равная 10^{-8} , приводит к уравнению $(x - 1)^4 = 10^{-8}$, имеющему следующие корни: $x_{1,2} = 1 \pm 10^{-2}$, $x_{3,4} = 1 \pm 10^{-2} \cdot i$. В этом случае погрешность, равная $10^{-6}\%$ в одном из коэффициентов, привела к погрешности решения примерно в 1 %, что явно говорит о плохой обусловленности задачи. ▲

3. Обусловленность задачи вычисления значения функции одной переменной. Пусть задача состоит в вычислении по заданному x значения $y = f(x)$ дифференцируемой функции f . В силу формул (2.21), (2.22) для этой задачи имеем:

$$v_{\Delta} \approx |f'(x)|, \quad (3.7)$$

$$v_{\delta} \approx \frac{|x| \cdot |f'(x)|}{|f(x)|}. \quad (3.8)$$

Воспользуемся этими формулами для оценки обусловленности задачи вычисления значений некоторых элементарных функций.

Пример 3.10. Для задачи вычисления значения функции $y = e^x$ в силу формулы (3.8) относительное число обусловленности v_{δ} приближенно равно $|x|$ и при реальных вычислениях оно не может быть очень большим. Например, при вычислении экспоненты на компьютере типа IBM PC всегда $|x| \lesssim 88$, так как в противном случае возможно переполнение или антiperеполнение. Следовательно, задача вычисления этой функции хорошо обусловлена, однако в случае $10 < |x| < 10^2$ следует ожидать потери одной-двух верных значащих цифр по сравнению с числом верных цифр аргумента x . Подчеркнем, что эта потеря точности объективно обусловлена погрешностью задания аргумента и не связана с используемым алгоритмом. ▲

Пример 3.11. Для задачи вычисления значения функции $y = \sin x$ в силу формулы (3.7) имеем $v_{\Delta} = |\cos x| \leq 1$, что говорит о хорошей абсолютной обусловленности этой задачи при всех x . Однако если важен результат с определенным числом верных знаков, то нужно исследовать относительную обусловленность. Согласно формуле (3.8) имеем $v_{\delta} \approx |x \operatorname{ctg} x|$. На рис. 3.3 приведен график этой функции при $x \geq 0$ (она четная).

Так как $v_{\delta} \rightarrow \infty$ при $x \rightarrow \pi k$ (для $k = \pm 1, \pm 2, \dots$), то при $x \approx \pi k$ задача обладает плохой относительной обусловленностью, и можно лишь утешаться тем, что мала абсолютная погрешность значения $y^* = \sin x^*$. Если же значение $|x|$ очень велико, то $v_{\delta} \gg 1$ и вычислять значение синуса просто бессмысленно. Например, если вычисления ведутся на компьютере типа IBM PC, где $\epsilon_m \sim 10^{-7}$ при вычислениях с одинарной точностью, то уже для $|x| \sim 10^7$ одна только абсолютная погрешность представления числа x есть величина порядка единицы, так как $\Delta(x^*) \approx |x| \cdot \epsilon_m$. Нетрудно понять, что при $\Delta(x^*) \sim 1$ вычисленное любым способом значение $\sin x^*$ не представляет никакой ценности. Вывод,

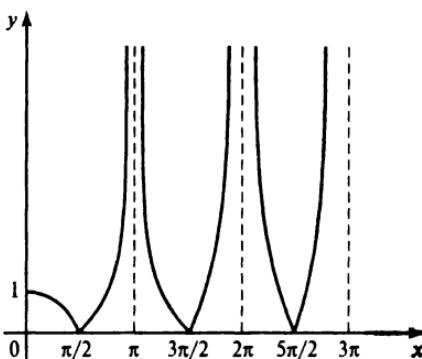


Рис. 3.3

который можно сделать, довольно прост: при использовании функции $y = \sin x$ желательно проводить вычисления так, чтобы аргумент находился в диапазоне $|x| < 2$, поскольку здесь $v_\delta \leq 1$. ▲

4. Обусловленность задачи вычисления интеграла $I = \int_a^b f(x) dx$.

Как следует из оценки (3.3), в этом случае абсолютное число обусловленности имеет вид $v_\Delta = b - a$. Если же перейти к рассмотрению относительных погрешностей и положить $\delta(f^*) = \sup_{[a, b]} |f^*(x) - f(x)| / |f(x)|$

для тех, x , где $f(x) \neq 0$, то используя неравенство

$$\Delta(I^*) \leq \int_a^b |f^*(x) - f(x)| dx \leq \int_a^b |f(x)| dx \cdot \delta(f^*),$$

получаем оценку

$$\delta(I^*) \leq v_\delta \delta(f^*), \quad (3.9)$$

$$\text{в которой } v_\delta = \int_a^b |f(x)| dx / \left| \int_a^b f(x) dx \right|.$$

Если подынтегральная функция знакопостоянна, то $v_\delta = 1$ и задача хорошо обусловлена. Если же функция f на отрезке $[a, b]$ принимает значения разных знаков, то $v_\delta > 1$. Для некоторых сильно колеблющихся (осциллирующих) около нуля функций может оказаться, что $v_\delta \gg 1$ и тогда задача вычисления интеграла является плохо обусловленной.

Иногда причиной появления плохо обусловленных задач становится отсутствие у пользователя компьютера элементарного представления об их существовании. В связи с этим заметим, что в последнее время получила развитие опасная тенденция пренебрежительного отношения к математическим знаниям вообще и к знанию вычислительной математики, в частности. Вера в могущество компьютеров и надежность стандартных программ бывает так велика, что считается совсем ненужным знать о вычисляемой математической величине что-либо кроме ее определения и, возможно, геометрического или физического смысла. Приведем простой пример, иллюстрирующий оценку (3.9) и заодно показывающий, что иногда аналитическая выкладка бывает эффективнее, чем применение самого современного компьютера.

Пример 3.12. Пусть в распоряжении пользователя имеется высокопроизводительный компьютер и стандартная программа для вычисления интегралов вида $I = \int_a^b f(x) dx$. Предположим, что нужно вычис-

лить коэффициенты Фурье $g_n = \int_{-1}^1 g(x) \sin(\pi n x) dx$ для функции $g(x) = x + A e^{\cos^2 x}$ при $n \gg 1$ и $A \gg 1$. Использование для вычисления g_n указанной стандартной программы при $f(x) = g(x) \sin(\pi n x)$ автоматически означает, что ставится вычислительная задача, рассмотренная выше. При этом относительная погрешность $\delta(f^*)$ вычисляемой на компьютере функции $f^* \approx f$ заведомо не может быть меньше¹ ε_m .

Оценим величину v_δ . Заметим, что функция $A e^{\cos^2 x} \sin(\pi n x)$ нечетная и, следовательно, интеграл от нее равен нулю. Поэтому возможно аналитическое вычисление интеграла:

$$g_n = \int_{-1}^1 x \sin(\pi n x) dx = -\frac{x}{\pi n} \cos(\pi n x) \Big|_{-1}^1 + \frac{1}{\pi n} \int_{-1}^1 \cos(\pi n x) dx = \frac{2}{\pi n} (-1)^{n+1}.$$

Кроме того, $|g(x)| \geq |A e^{\cos^2 x}| - |x| \geq A - 1$ и поэтому

$$\int_{-1}^1 |f(x)| dx = \int_{-1}^1 |g(x)| \cdot |\sin(\pi n x)| dx \geq (A - 1) \int_{-1}^1 |\sin(\pi n x)| dx = A - 1.$$

¹ Если учитывать погрешность вычисления функции $\sin(\pi n x)$, то следует ожидать, что эта погрешность будет значительно больше ε_m .

Таким образом, $v_\delta \geq (A - 1)\pi n/2 \approx 1.5An$. Если g_n вычисляется, например при $A = 10^4$ и $n = 10^2$, то $v_\delta \geq 1.5 \cdot 10^6$. Следовательно, для этих значений параметров принятное решение о вычислении интеграла по простейшей стандартной программе может обойтись в потерю примерно шести значащих цифр результата. Если вычисления ведутся на компьютере, имеющем лишь около семи десятичных значащих цифр, то возможна катастрофическая потеря точности. Заметим, что в этих рассуждениях никак не учтены погрешность реализованного в программе метода и вычислительная погрешность, наличие которых усугубляет ситуацию. ▲

5. Обусловленность задачи вычисления суммы ряда. Рассмотрим задачу вычисления суммы абсолютно сходящегося ряда с ненулевыми слагаемыми. В силу оценки (3.4) эта задача устойчива, а за относительное число обусловленности следует принять величину

$$v_\delta = \sum_{k=0}^{\infty} |a_k| / \left| \sum_{k=0}^{\infty} a_k \right|. \quad (3.10)$$

Заметим, что для ряда с положительными слагаемыми имеем $v_\delta = 1$, т.е. задача хорошо обусловлена. Если же ряд незнакопостоянный, то $v_\delta > 1$ и при $v_\delta \gg 1$ задача оказывается плохо обусловленной.

§ 3.3. Вычислительные методы

Обсудив некоторые важные особенности вычислительных задач, обратим внимание на те методы, которые используются в вычислительной математике для преобразования задач к виду, удобному для реализации на компьютере, и позволяют конструировать вычислительные алгоритмы. Мы будем называть эти методы *вычислительными*. С некоторой степенью условности можно разбить вычислительные методы на следующие классы:

- 1) методы эквивалентных преобразований;
- 2) методы аппроксимации;
- 3) прямые (точные) методы;
- 4) итерационные методы;
- 5) методы статистических испытаний (методы Монте-Карло).

Метод, осуществляющий вычисление решения конкретной задачи, может иметь довольно сложную структуру, но его элементарными шагами, являются, как правило, реализации указанных методов. Дадим о них общее представление.

1. Методы эквивалентных преобразований. Эти методы позволяют заменить исходную задачу другой, имеющей то же решение.

Выполнение эквивалентных преобразований оказывается полезным, если новая задача проще исходной или обладает лучшими свойствами, или для нее существует известный метод решения, а, может быть, и готовая программа.

Пример 3.13. Эквивалентное преобразование квадратного уравнения

$$x^2 + bx + c = 0 \text{ к виду } \left(x + \frac{b}{2}\right)^2 = \frac{b^2 - 4c}{4} \text{ (выделение полного квадрата)}$$

сводит задачу к проблеме вычисления квадратного корня и приводит к известным для ее корней формулам (3.2). ▲

Эквивалентные преобразования иногда позволяют свести решение исходной вычислительной задачи к решению вычислительной задачи совершенно иного типа.

Пример 3.14. Задача отыскания корня нелинейного уравнения $f(x) = 0$ может быть сведена к эквивалентной задаче поиска точки глобального минимума функции $\Phi(x) = (f(x))^2$. В самом деле, функция $\Phi(x)$ неотрицательна и достигает минимального значения, равного нулю, при тех и только тех x , для которых $f(x) = 0$. ▲

2. Методы аппроксимации. Эти методы позволяют приблизить (аппроксимировать) исходную задачу другой, решение которой в определенном смысле близко к решению исходной задачи. Погрешность, возникающая при такой замене, называется *погрешностью аппроксимации*. Как правило, аппроксимирующая задача содержит некоторые параметры, позволяющие регулировать значение погрешности аппроксимации или воздействовать на другие свойства задачи. Принято говорить, что метод аппроксимации *сходится*, если погрешность аппроксимации стремится к нулю при стремлении параметров метода к некоторому предельному значению.

Пример 3.15. Один из простейших способов вычисления интеграла $I = \int_a^b f(x) dx$ состоит в аппроксимации интеграла на основании формулы прямоугольников величиной

$$I^h = h \sum_{i=1}^n f\left(a + \left(i - \frac{1}{2}\right) h\right).$$

Шаг $h = (b - a)/n$ является здесь параметром метода. Так как I^h представляет собой специальным образом построенную интегральную

сумму, то из определения определенного интеграла следует, что $I^h \rightarrow I$ при $h \rightarrow 0$, т.е. метод прямоугольников сходится. ▲

Пример 3.16. Учитывая определение производной функции $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$, для ее приближенного вычисления можно использовать формулу $f'(x) \approx \frac{f(x+h) - f(x)}{h}$. Погрешность аппроксимации этой формулы численного дифференцирования стремится к нулю при $h \rightarrow 0$. ▲

Одним из распространенных методов аппроксимации является *дискретизация* — приближенная замена исходной задачи *конечномерной задачей*, т.е. задачей, входные данные и искомое решение которой могут быть однозначно заданы конечным набором чисел. Для задач, которые не являются конечномерными, этот шаг необходим для последующей реализации на компьютере, так как вычислительная машина в состоянии оперировать лишь с конечным количеством чисел. В приведенных выше примерах 3.15 и 3.16 была использована дискретизация. Хотя точное вычисление интеграла и предполагает использование бесконечного множества значений $f(x)$ (для всех $x \in [a, b]$), его приближенное значение можно вычислить, использовав конечное число n значений в точках $a + (i - 1/2)h$. Аналогично, задача вычисления производной, точное решение которой предполагает выполнение операции предельного перехода при $h \rightarrow 0$ (а следовательно, использование бесконечного множества значений функции f), сводится к приближенному вычислению производной по двум значениям функции.

При решении нелинейных задач широко используют различные *методы линеаризации*, состоящие в приближенной замене исходной задачи более простыми линейными задачами.

Пример 3.17. Пусть требуется приближенно вычислить значение $x = \sqrt{a}$ для $a > 0$ на компьютере, способном выполнять простейшие арифметические операции. Заметим, что по определению $x = \sqrt{a}$ является положительным корнем нелинейного уравнения $x^2 - a = 0$. Пусть $x^{(0)}$ — некоторое известное приближение к \sqrt{a} . Заменим параболу $y = x^2 - a$ прямой $y = (x^{(0)})^2 - a + 2x^{(0)}(x - x^{(0)})$, являющейся касательной, проведенной к ней в точке с абсциссой $x = x^{(0)}$ (рис. 3.4). Точка пересечения этой касательной с осью Ox дает лучшее, чем $x^{(0)}$, приближение и находится из

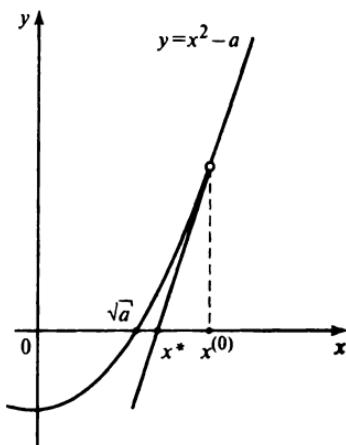


Рис. 3.4

линейного уравнения $(x^{(0)})^2 - a + 2x^{(0)}(x - x^{(0)}) = 0$. Решая его, получаем приближенную формулу

$$\sqrt{a} \approx x^* = x^{(0)} - \frac{(x^{(0)})^2 - a}{2x^{(0)}} = \frac{1}{2} \left(x^{(0)} + \frac{a}{x^{(0)}} \right). \quad (3.11)$$

Например, если для $x = \sqrt{2}$ взять $x^{(0)} = 2$, то получится уточненное значение $\sqrt{2} \approx x^* = \frac{1}{2} \left(2 + \frac{2}{2} \right) = 1.5$. ▲

При решении разных классов вычислительных задач могут использоваться различные методы аппроксимации; к ним можно отнести и методы регуляризации решения некорректных задач. Заметим, что методы регуляризации широко используются для решения плохо обусловленных задач.

3. Прямые методы. Метод решения задачи называют *прямым*, если он позволяет получить решение после выполнения конечного числа элементарных операций

Пример 3.18. Метод вычисления корней квадратного уравнения $x^2 + bx + c = 0$ по формулам $x_{1,2} = (-b \pm \sqrt{b^2 - 4c})/2$ является прямым методом. Элементарными здесь считаются четыре арифметические операции и операция извлечения квадратного корня. ▲

Заметим, что элементарная операция прямого метода может оказаться довольно сложной (вычисление значений элементарной или специальной функции, решение системы линейных алгебраических урав-

нений, вычисление определенного интеграла и т.д.). То, что она принимается за элементарную, предполагает во всяком случае, что ее выполнение существенно проще вычисления решения всей задачи.

При построении прямых методов существенное внимание уделяется минимизации числа элементарных операций.

Пример 3.19. (схема Горнера¹). Пусть задача состоит в вычислении значения многочлена

$$P_n(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n \quad (3.12)$$

по заданным коэффициентам a_0, a_1, \dots, a_n и значению аргумента x . Если вычислять многочлен непосредственно по формуле (3.12), причем x^2, x^3, \dots, x^n находить последовательным умножением на x , то потребуется выполнить $2n - 1$ операций умножения и n операций сложения.

Значительно более экономичным является метод вычисления, называемый *схемой Горнера*. Он основан на записи многочлена в следующем эквивалентном виде:

$$P_n(x) = ((\dots ((a_0x + a_{n-1})x + a_{n-2})x + \dots)x + a_1)x + a_0.$$

Расстановка скобок диктует такой порядок вычислений: $S_0 = a_n, S_1 = S_0x + a_{n-1}, S_2 = S_1x + a_{n-2}, \dots, S_i = S_{i-1}x + a_{n-i}, \dots, S_n = S_{n-1}x + a_0$. Здесь вычисление значения $P_n(x) = S_n$ потребовало выполнения только n операций умножения и n операций сложения.

Схема Горнера интересна тем, что дает пример оптимального по числу элементарных операций метода. В общем случае значение $P_n(x)$ нельзя получить никаким методом в результате выполнения меньшего числа операций умножения и сложения. ▲

Иногда прямые методы называют *точными*, подразумевая под этим, что при отсутствии погрешностей во входных данных и при точном выполнении элементарных операций полученный результат также будет точным. Однако при реализации метода на компьютере неизбежно появление вычислительной погрешности, которая зависит от чувствительности метода к погрешностям округления. Многие прямые (точные) методы, разработанные в домашинный период, оказались непригодными для машинных вычислений именно из-за чрезмерной чувствительности к погрешностям округления. Не все точные методы таковы, однако стоит заметить, что не совсем удачный термин «точный» характеризует свойства идеальной реализации метода, но отнюдь не качество полученного при реальных вычислениях результата.

¹ Вильямс Джордж Горнер (1786—1837) — английский математик.

4. Итерационные методы. Это — специальные методы построения последовательных приближений к решению задачи. Применение метода начинают с выбора одного или нескольких начальных приближений. Для получения каждого из последующих приближений выполняют однотипный набор действий с использованием найденных ранее приближений — *итерацию*¹. Неограниченное продолжение этого *итерационного процесса* теоретически позволяет построить бесконечную последовательность приближений к решению — *итерационную последовательность*. Если эта последовательность сходится к решению задачи, то говорят, что *итерационный метод сходится*. Множество начальных приближений, для которых метод сходится, называется *областью сходимости метода*.

Заметим, что итерационные методы широко используются при решении самых разнообразных задач с применением компьютеров.

Пример 3.20. Рассмотрим известный итерационный метод, предназначенный для вычисления \sqrt{a} (где $a > 0$), — *метод Ньютона*. Зададим произвольное начальное приближение $x^{(0)} > 0$. Следующее приближение вычислим по формуле $x^{(1)} = \frac{1}{2} \left(x^{(0)} + \frac{a}{x^{(0)}} \right)$, выведенной с помощью метода линеаризации в примере 3.17 (см. формулу (3.11)). Продолжив этот процесс далее, получим итерационную последовательность $x^{(0)}, x^{(1)}, x^{(2)}, \dots, x^{(n)}, \dots$, в которой очередное $(k+1)$ -е приближение вычисляется через k -е по *рекуррентной*² формуле

$$x^{(k+1)} = \frac{1}{2} \left(x^{(k)} + \frac{a}{x^{(k)}} \right), \quad k \geq 0. \quad (3.13)$$

Известно, что этот метод сходится при любом начальном приближении $x^{(0)} > 0$, так что его область сходимости — множество всех положительных чисел.

Вычислим с его помощью значение $\sqrt{2}$ на 8-разрядном десятичном компьютере. Зададим $x^{(0)} = 2$ (как в примере 3.17). Тогда $x^{(1)} = 1.5$, $x^{(2)} = 1.4166667$, $x^{(3)} = 1.4142157$, $x^{(4)} = 1.4142136$, $x^{(5)} = 1.4142136$. Дальнейшие вычисления бессмысленны, так как из-за ограниченности разрядной сетки все следующие уточнения будут давать тот же результат. Однако сравнение с точным значением $\sqrt{2} = 1.41421356\dots$ показывает, что уже на третьей итерации были получены шесть верных значащих цифр. ▲

¹ От латинского слова *iteratio* — «повторение».

² От латинского слова *recurrent* — «возвращающийся».

Обсудим на примере метода Ньютона некоторые типичные для итерационных методов (и не только для них) проблемы. Итерационные методы по своей сути являются приближенными; ни одно из получаемых приближений не является точным значением решения. Однако сходящийся итерационный метод дает принципиальную возможность найти решение с любой заданной точностью $\epsilon > 0$. Поэтому, применяя итерационный метод, всегда задают требуемую точность ϵ и итерационный процесс прерывают, как только она достигается.

Хотя сам факт сходимости метода безусловно важен, он недостаточен для того, чтобы рекомендовать метод для использования на практике. Если метод сходится очень медленно (например, для получения решения с точностью в 1 % нужно сделать 10^6 итераций), то он непригоден для вычислений на компьютере. Практическую ценность представляют быстро сходящиеся методы, к которым относится и метод Ньютона (напомним, что в рассмотренном выше примере точность $\epsilon = 10^{-5}$ в вычислении $\sqrt{2}$ была достигнута всего за три итерации). Для теоретического исследования скорости сходимости и условий применимости итерационных методов выводят так называемые *априорные¹ оценки погрешности*, позволяющие еще до вычислений дать некоторое заключение о качестве метода.

Приведем две такие априорные оценки для метода Ньютона. Пусть $x^{(0)} > \sqrt{a}$. Известно, что тогда $x^{(n)} > \sqrt{a}$ для всех $n \geq 0$ и погрешности двух последовательных приближений связаны следующим неравенством:

$$\delta^{(n)} \leq (\delta^{(n-1)})^2. \quad (3.14)$$

Здесь $\delta^{(n)} = \frac{1}{2\sqrt{a}} (x^{(n)} - \sqrt{a}) = \frac{\delta(x^{(n)})}{2}$ — величина, характеризующая

относительную погрешность n -го приближения. Это неравенство говорит об очень высокой квадратичной скорости сходимости метода: на каждой итерации «ошибка» $\delta^{(n)}$ возводится в квадрат. Если выразить $\delta^{(n)}$ через погрешность начального приближения, то получим неравенство

$$\delta^{(n)} \leq (\delta^{(0)})^{2^n}, \quad (3.15)$$

из которого видна роль хорошего выбора начального приближения. Чем меньше $\delta^{(0)} = \delta(x^{(0)})/2$, тем быстрее будет сходиться метод.

Практическая реализация итерационных методов всегда связана с необходимостью выбора *критерия окончания итерационного процесса*. Вычисления не могут продолжаться бесконечно долго и должны быть прерваны в соответствии с некоторым критерием, связанным, например с достижением заданной точности. Использование для этой

¹ От латинского слова *a'priori* — «до опыта».

цели априорных оценок чаще всего оказывается невозможным или неэффективным. Качественно верно описывая поведение метода, такие оценки являются завышенными и дают весьма недостоверную количественную информацию. Нередко априорные оценки содержат неизвестные величины (например, в оценках (3.14), (3.15) содержится величина \sqrt{a}), либо предполагают наличие и серьезное использование некоторой дополнительной информации о решении. Чаще всего такой информации нет, а ее получение связано с необходимостью решения дополнительных задач, нередко более сложных, чем исходная.

Для формирования критерия окончания по достижении заданной точности, как правило, используют так называемые *апостериорные¹ оценки погрешности* — неравенства, в которых значение погрешности оценивается через известные или получаемые в ходе вычислительного процесса величины. Хотя такими оценками нельзя воспользоваться до начала вычислений, в ходе вычислительного процесса они позволяют давать конкретную количественную оценку погрешности.

Например, для метода Ньютона (3.13) справедлива следующая апостериорная оценка:

$$|x^{(n)} - \sqrt{a}| \leq |x^{(n)} - x^{(n-1)}|, \quad n \geq 1,$$

позволяющая оценивать абсолютную погрешность приближения через модуль разности двух последовательных приближений. Она дает возможность сформулировать при заданной точности $\varepsilon > 0$ очень простой критерий окончания. Как только окажется выполненным неравенство

$$|x^{(n)} - x^{(n-1)}| < \varepsilon, \quad (3.16)$$

вычисления следует прекратить и принять $x^{(n)}$ за приближение к \sqrt{a} с точностью ε .

Пример 3.21. Если значение $\sqrt{2}$ требуется найти с точностью $\varepsilon = 10^{-5}$, то неравенство (3.16), как видно из приведенных в примере 3.20 вычислений, будет выполнено при $n = 4$. Так как $|x^{(4)} - x^{(3)}| = 3 \cdot 10^{-6}$, то с учетом погрешности округления можно записать результат: $\sqrt{2} = 1.41421 \pm 0.00001$. ▲

5. Методы статистических испытаний (методы Монте-Карло). Это — численные методы, основанные на моделировании случайных величин и построении статистических оценок решений задач. Этот класс

¹ От латинского слова *a'priori* — «после опыта».

методов, как принято считать, возник в 1949 г., когда Дж. фон Нейман¹ и С. Уlam² использовали случайные числа для моделирования с помощью компьютеров поведения нейтронов в ядерном реакторе. Эти методы могут оказаться незаменимыми при моделировании больших систем, но подробное их изложение предполагает существенное использование аппарата теории вероятностей и математической статистики и выходит за рамки данной книги.

§ 3.4. Корректность вычислительных алгоритмов

1. Вычислительный алгоритм. Вычислительный метод, доведенный до степени детализации, позволяющей реализовать его на компьютере, принимает форму вычислительного алгоритма.

Определим *вычислительный алгоритм* как точное предписание действий над входными данными, задающее вычислительный процесс, направленный на преобразование произвольных входных данных x (из множества допустимых для данного алгоритма входных данных X) в полностью определяемый этими входными данными результат.

Реальный вычислительный алгоритм складывается из двух частей: *абстрактного вычислительного алгоритма*, формулируемого в общепринятых математических терминах, и *программы*, записанной на одном из алгоритмических языков и предназначенной для реализации алгоритма на компьютере. Как правило, в руководствах по методам вычислений излагаются именно абстрактные алгоритмы, но их обсуждение проводится так, чтобы выявить особенности алгоритмов, которые оказывают существенное влияние на качество программной реализации.

2. Определение корректности алгоритма. К вычислительным алгоритмам, предназначенным для широкого использования, предъявляется ряд весьма жестких требований. Первое из них — корректность алгоритма. Будем называть вычислительный алгоритм *корректным*, если выполнены три условия:

1) он позволяет после выполнения конечного числа элементарных для вычислительной машины операций преобразовать любое входное данное $x \in X$ в результат y ;

¹ Джон фон Нейман (1903—1957) — американский математик, физик, инженер-изобретатель. Оказал значительное влияние на развитие современной математики. Один из создателей первых компьютеров. Историческую справку о фон Неймане см. в [54].

² Станислав Улам (1909—1984) — американский математик. Краткий очерк о его деятельности можно найти в [54].

2) результат устойчив по отношению к малым возмущениям входных данных;

3) результат обладает вычислительной устойчивостью.

Если хотя бы одно из перечисленных условий не выполнено, то будем называть алгоритм *некорректным*. Уточним и более подробно обсудим эти условия.

Необходимость выполнения первого условия понятна. Если для получения результата нужно выполнить бесконечное число операций либо требуются операции, не реализованные на компьютерах, то алгоритм следует признать некорректным.

Пример 3.22. Известный алгоритм деления чисел «углом» некорректен, так как он может продолжаться бесконечно, если не определен критерий окончания вычислений. ▲

Пример 3.23. Отсутствие критерия окончания делает некорректным и алгоритм Ньютона вычисления \sqrt{a} (см. пример 3.20). ▲

Пример 3.24. Алгоритм вычисления корней квадратного уравнения (3.1) по формулам (3.2) некорректен, если он предназначен для использования на вычислительной машине, на которой не реализована операция извлечения квадратного корня. ▲

3. Устойчивость по входным данным. Устойчивость результата у к малым возмущениям входных данных (*устойчивость по входным данным*) означает, что результат непрерывным образом зависит от входных данных при условии, что отсутствует вычислительная погрешность. Это требование устойчивости аналогично требованию устойчивости вычислительной задачи. Отсутствие такой устойчивости делает алгоритм непригодным для использования на практике.

Отметим, что в формулировку устойчивости алгоритма по входным данным неявно входит одно весьма важное предположение, а именно, что вместе с входным данным x в множество допустимых входных данных X входят и все близкие к x приближенные входные данные x^* .

Пример 3.25. Пусть алгоритм предназначен для вычисления корней квадратного уравнения (3.1) с коэффициентами, удовлетворяющими условию $d = b^2 - 4c \geq 0$. Если в нем используются формулы (3.2), то этот алгоритм некорректен. В самом деле, значение d^* , отвечающее приближенно заданным коэффициентам b^* и c^* , может оказаться отрицательным, если $d \approx 0$. Тогда вычисления завершатся аварийным остановом при попытке извлечь квадратный корень из отрицательного числа. Если же в формуле (3.2) заменить d на $\max\{d, 0\}$, то алгоритм становится корректным. ▲

4. Вычислительная устойчивость. Из-за наличия погрешностей округления при вводе входных данных в компьютер и при выполнении арифметических операций неизбежно появление вычислительной погрешности. На разных компьютерах она различна из-за различий в разрядности и способах округления, но для фиксированного алгоритма в основном значение погрешности определяется машинной точностью ϵ_m .

Назовем алгоритм *вычислительно устойчивым*, если вычислительная погрешность результата стремится к нулю при $\epsilon_m \rightarrow 0$. Обычно вычислительный алгоритм называют *устойчивым*, если он устойчив по входным данным и вычислительно устойчив, и — *неустойчивым*, если хотя бы одно из этих условий не выполнено.

Пример 3.26.¹ Пусть требуется составить таблицу значений интегралов $I_n = \int_0^1 x^n e^{1-x} dx$ для $n = 1, 2, \dots$ на 6-разрядном десятичном компьютере.

Интегрируя по частям, имеем

$$I_n = \int_0^1 x^n d(-e^{1-x}) = -x^n e^{1-x} \Big|_0^1 + \int_0^1 e^{1-x} d(x^n) = -1 + \int_0^1 n x^{n-1} e^{1-x} dx.$$

Следовательно, справедлива формула

$$I_n = n I_{n-1} - 1, \quad n \geq 1. \quad (3.17)$$

Кроме того,

$$I_0 = \int_0^1 e^{1-x} dx = e - 1 \approx I_0^* = 1.71828.$$

Воспользуемся формулой (3.17) для последовательного вычисления приближенных значений интегралов I_n :

$$\begin{aligned} I_1 &\approx I_1^* = 1I_0^* - 1 = 0.71828; & I_2 &\approx I_2^* = 2I_1^* - 1 = 0.43656; \\ I_3 &\approx I_3^* = 3I_2^* - 1 = 0.30968; & I_4 &\approx I_4^* = 4I_3^* - 1 = 0.23872; \\ I_5 &\approx I_5^* = 5I_4^* - 1 = 0.19360; & I_6 &\approx I_6^* = 6I_5^* - 1 = 0.16160; \\ I_7 &\approx I_7^* = 7I_6^* - 1 = 0.13120; & I_8 &\approx I_8^* = 8I_7^* - 1 = 0.04960; \\ I_9 &\approx I_9^* = 9I_8^* - 1 = -0.55360; & I_{10} &\approx I_{10}^* = 10I_9^* - 1 = -6.5360. \end{aligned}$$

¹ Идея примера заимствована из [5, 106].

Здесь вычисления следует прекратить. Искомые значения интегралов очевидно, положительны, а найденные значения при $n = 9$ и $n = 10$ отрицательны. В чем причина появления такой большой погрешности?

В данном примере все вычисления проводились точно, а единственная и, на первый взгляд, незначительная ошибка была сделана при округлении значения I_0 до шести значащих цифр (заметим, что $\Delta_0 = |I_0 - I_0^*| \approx 2 \cdot 10^{-6}$). Однако при вычислении I_1 эта погрешность сохранилась, при вычислении I_2 умножилась на $2!$, при вычислении I_3 — на $3!$, ..., при вычислении I_9 — на $9!$ и т.д.

Таким образом, $\Delta_n = |I_n - I_n^*| = n! \Delta_0$. Уже при $n = 9$ имеем $9! = 36\,880$ и поэтому $\Delta_9 = 9! \Delta_0 \approx 0.73$.

Если вычисления производятся без ограничений на число n , то рассматриваемый алгоритм следует признать вычислительно неустойчивым. Погрешности растут пропорционально $n!$ настолько быстро, что уже при довольно скромных значениях n попытки добиться приемлемого результата даже за счет увеличения разрядности мантиссы заранее обречены на неудачу,

Как изменить алгоритм, чтобы сделать его устойчивым? Перепишем формулу (3.17) в виде

$$I_{n-1} = \frac{I_n + 1}{n}, \quad n \geq 1 \quad (3.18)$$

и будем вести вычисления значений I_n в обратном порядке, начиная,

например с $n = 54$. Положим $I_{54} \approx I_{54}^* = 0$. Так как $I_{54} \leq e \int_0^1 x^{54} dx = e/55$, то $\Delta_{54} \leq e/55 \approx 5 \cdot 10^{-2}$. Однако при вычислении I_{53} эта погрешность уменьшится в 54 раза, при вычислении I_{52} — еще в 53 раза и т.д.

В результате значения I_n при $n = 50, \dots, 1$ будут вычислены с шестью верными значащими цифрами. Здесь погрешности не растут, а затухают. Ясно, что модифицированный алгоритм вычислительно устойчив. ▲

Вычислительная неустойчивость алгоритма часто может быть выявлена благодаря анализу устойчивости по входным данным, так как неустойчивость к малым погрешностям округления входных данных автоматически свидетельствует о вычислительной неустойчивости алгоритма.

Пример 3.27. Предположим, что значения y_n для $n = 1, 2, \dots$ вычисляются по рекуррентной формуле

$$y_n = \alpha_n y_{n-1} + \beta_n, \quad (3.19)$$

а величина y_0 задана. Пусть y_0^* — заданное приближенное значение величины y_0 . Тогда (если вычисления ведутся абсолютно точно) определяемые по формуле (3.19) приближенные значения содержат погрешности, связанные равенством $y_n - y_n^* = \alpha_n(y_{n-1} - y_{n-1}^*)$. Следовательно, $\Delta(y_n^*) = |\alpha_n| \Delta(y_{n-1}^*)$ и при выполнении условия $|\alpha_n| \leq 1$ алгоритм устойчив по входным данным, поскольку $\Delta(y_n^*) \leq \Delta(y_0^*)$ для всех n . Если же $|\alpha_n| \geq q > 1$, то $\Delta(y_n^*) \geq q^n \Delta(y_0^*)$ и абсолютная погрешность неограниченно возрастает при $n \rightarrow \infty$. В этом случае алгоритм неустойчив по входным данным, а потому и вычислительно неустойчив. ▲

Справедливости ради следует заметить, что алгоритм (3.19) был признан нами неустойчивым в случае $|\alpha_n| \geq q > 1$ при выполнении двух условий, на которых не было достаточно акцентировано внимание. Первое из них состоит в предположении о неограниченной продолжительности вычислительного процесса ($n \rightarrow \infty$), что невозможно на практике. В действительности такой характер неустойчивости говорит о тенденции к неограниченному росту погрешности при неограниченном продолжении вычислений. Второе условие касается выбранной меры погрешности. Совсем не обязательно, чтобы рост абсолютной погрешности всегда был неприемлем в конкретных вычислениях. Если он сопровождается сильным ростом точного решения и при этом относительная погрешность остается малой, то алгоритм можно признать относительно устойчивым. По-видимому, при анализе вычислительной устойчивости более естественным является рассмотрение относительных погрешностей.

Пример 3.28. Пусть в формуле (3.19) все $\beta_n = 0$ и $\alpha_n \neq 0$. Тогда $\delta(y_n^*) = \Delta(y_n^*) / |y_n| = |\alpha_n| \Delta(y_{n-1}^*) / |\alpha_n y_{n-1}| = \delta(y_{n-1}^*)$. Следовательно, $\delta(y_n^*) = \delta(y_0^*)$ и при любых значениях $\alpha_n \neq 0$ алгоритм относительно устойчив по входным данным. ▲

§ 3.5. Чувствительность вычислительных алгоритмов к ошибкам округления

Выполнение вычислений на компьютере сопровождается появлением вычислительной погрешности, связанной в первую очередь с необходимостью округления результата каждой арифметической операции. Даже если разрядность компьютера велика, существует реальная опасность, что выполнение большого числа операций приведет к накоплению погрешностей, способных значительно или даже полностью иска-

зить вычисляемый результат. Однако и при небольшом числе действий результат вычислений может оказаться совершенно неправильным, если алгоритм слишком чувствителен к ошибкам округления.

Начинающий вычислитель часто склонен игнорировать ошибки округления. На первых порах при решении простых задач, в особенности если новичок не задумывается о точности найденных решений, его позицию нетрудно понять. Однако решение серьезных задач (когда число арифметических операций превышает миллиарды, в вычисления вкладываются значительные средства, а результат следует получить с гарантированной точностью и за принятые на основании расчетов решения приходится нести ответственность) предполагает совсем иное отношение к вычислительной погрешности.

1. Порядок выполнения операций. Решение математической задачи на компьютере сводится в конечном итоге к выполнению последовательности простейших арифметических и логических операций. Однако часто одно и то же математическое выражение допускает различные способы вычисления, отличающиеся только порядком выполнения операций. Если вычисления производить точно, то они (при любом способе вычисления) будут приводить к одному результату. Однако результаты вычислений на компьютере уже зависят от порядка выполнения операций и различие в вычислительной погрешности может быть весьма значительным.

Рассмотрим простой, но полезный пример вычисления на компьютере суммы $S_N = \sum_{i=1}^N a_i$. Пусть a_i — положительные представимые на вычислительной машине числа. В каком порядке следует их суммировать для того, чтобы сделать вычислительную погрешность по возможности минимальной?

Пусть $S_k = \sum_{i=1}^k a_i$ — частичная сумма, а S_k^* — ее приближенное значение, вычисляемое по формуле $S_k^* = S_{k-1}^* \otimes a_k$. Погрешность значения S_k^* складывается из погрешности значения S_{k-1}^* и погрешности выполнения операции $S_{k-1}^* \otimes a_k$. Следовательно, $\bar{\Delta}(S_k^*) = \bar{\Delta}(S_{k-1}^*) + (S_{k-1}^* + a_k)\varepsilon_m \approx \bar{\Delta}(S_{k-1}^*) + S_k\varepsilon_m$. Поэтому $\Delta(S_N^*) \approx (S_N\varepsilon_m + S_{N-1}\varepsilon_m + \dots + S_2\varepsilon_m) = ((N-1)a_1 + (N-1)a_2 + (N-2)a_3 + \dots + 2a_{N-1} + a_N)\varepsilon_m$.

Так как множитель, с которым входит a_i в формулу для оценки погрешности, убывает с ростом i , в общем случае погрешность ока-

жется наименьшей, если суммировать числа в порядке возрастания их значений, начиная с наименьшего.

Иногда неудачно выбранный порядок операций либо приводит к полной потере точности, либо вообще не дает возможности получить результат из-за переполнения.

Пример 3.29. Пусть на компьютере типа IBM PC требуется вычислить произведение $v = a_0 \times a_1 \times a_2 \times \dots \times a_{49} \times a_{50}$, где $a_i \approx 10^{25-i}$. Если производить вычисления в естественном порядке, то уже $a_0 \times a_1 \approx 10^{44}$ даст значение, равное $+\infty$. Вычисление произведения в обратном порядке сразу же приводит к исчезновению порядка, так как $a_{50} \times a_{49} \approx 10^{-49} < X_0$. В данном случае вполне приемлем следующий порядок операций: $v = a_0 \times a_{50} \times a_1 \times a_{49} \times \dots \times a_{24} \times a_{26} \times a_{25}$, исключающий возможность переполнения или антипереполнения. ▲

2. Катастрофическая потеря точности. Иногда короткая последовательность вычислений приводит от исходных данных, известных с высокой точностью, к результату, содержащему недопустимо мало верных цифр или вообще не имеющему ни одной верной цифры. В этом случае, как было отмечено в § 3.2, принято говорить о катастрофической потере точности.

Пример 3.30¹. Известно, что функция e^x может быть представлена в виде сходящегося степенного ряда:

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!} + \dots \quad (3.20)$$

Возможность вычисления значения экспоненты прямым суммированием ряда (3.20) кажется привлекательной. Пусть для вычислений используется 6-разрядный десятичный компьютер. Возьмем $x = -8.1$ и будем вычислять значения частичных сумм до тех пор, пока добавление очередного слагаемого еще меняет значение суммы:

$$\begin{aligned} e^{-8.1} \approx 1.00000 &\ominus 8.10000 \oplus 32.8050 \ominus 88.5737 \oplus 179.362 \ominus 290.566 \oplus \dots \\ &\dots \oplus 16.4111 \ominus 7.81941 \oplus \dots = 0.000649915. \end{aligned}$$

В сумму вошло 36 слагаемых и значение очередного (37-го) слагаемого оказалось уже не в состоянии изменить результат. Можно ли считать результат удовлетворительным? Сравнение с истинным значением $e^{-8.1} \approx 0.000303539$ показывает, что найденное значение не содержит ни одной верной цифры.

¹ Идея примера заимствована из [106].

В чем причина катастрофической потери точности? Дело в том, что вычисленные слагаемые ряда неизбежно содержат погрешности, причем для некоторых (для слагаемых с 5-го по 13-е) погрешность превосходит значение самого искомого результата. Налицо явный дефект алгоритма, к обсуждению которого мы еще вернемся в конце этого параграфа.

В данном случае переход к вычислениям с удвоенной длиной мантиссы позволит получить значение $e^{-8.1}$ с шестью верными значащими цифрами. Однако всего лишь удвоенное значение аргумента $x = -16.2$ снова возвращает нас к той же проблеме. Поступим иначе. Используя разложение (3.20), вычисляем $e^{8.1} \approx 1.00000 \oplus 8.1000 \oplus 32.8050 \oplus \dots = 3294.47$ и тогда $e^{-8.1} = 1/e^{8.1} \approx 0.000303539$. Предложенное изменение алгоритма позволило получить искомое значение на том же 6-разрядном десятичном компьютере, но уже с шестью верными значащими цифрами.

Заметим тем не менее, что реальные машинные алгоритмы вычисления e^x устроены совсем иначе. ▲

Приведем теперь пример, когда к катастрофической потере точности приводит еще более короткая последовательность вычислений.

Пример 3.31. Пусть при $x = 1/490$ на 6-разрядном десятичном компьютере вычисляется значение функции

$$y = \cos x - \cos 2x. \quad (3.21)$$

Заметим, что при $x \approx 0$ величины $c_1 = \cos x \approx 1 - x^2/2$ и $c_2 = \cos 2x \approx 1 - 2x^2$ близки. Так как вычисленные их приближенные значения c_1^* и c_2^* будут содержать погрешность, то возможна серьезная потеря точности. Действительно, $c_1^* = 0.999998$, $c_2^* = 0.999992$ и $y^* = c_1^* - c_2^* = 0.000006$. При вычитании старшие разряды оказались потерянными и в результате осталась только одна значащая цифра.

Вычисление по эквивалентной формуле $y = 2 \sin(x/2) \sin(3x/2)$, позволяет избежать вычитания близких чисел и дает значение $y^* = 0.624741 \cdot 10^{-5}$ с шестью верными цифрами.

Интересно отметить, что $y \approx 1.5x^2$, причем использование этой приближенной формулы в данном случае дает шесть верных значащих цифр, в то время как вычисления по формуле (3.21) — только одну верную цифру. ▲

Замечание. Не всегда катастрофическая потеря точности в промежуточных вычислениях действительно является катастрофой. Все зависит от того, как в дальнейшем используется результат.

3. Обусловленность вычислительного алгоритма. По аналогии с понятием обусловленности математической задачи можно ввести понятие обусловленности вычислительного алгоритма, отражающее чувствительность результата работы алгоритма к малым, но неизбежным погрешностям округления. Вычислительно устойчивый алгоритм называют *хорошо обусловленным*, если малые относительные погрешности округления (характеризуемые числом ε_m) приводят к малой относительной вычислительной погрешности $\delta(y^*)$ результата y^* , и плохо обусловленным, если вычислительная погрешность может быть недопустимо большой.

Если $\delta(y^*)$ и ε_m связаны неравенством $\delta(y^*) \leq v_A \varepsilon_m$, то число v_A следует называть *числом обусловленности вычислительного алгоритма*. Для плохо обусловленного алгоритма $v_A \gg 1$.

При очень большом значении числа обусловленности алгоритм можно считать практически неустойчивым¹.

Применим, например алгоритм, первоначально предложенный в примере 3.26, для вычисления конечной серии из N интегралов I_1, I_2, \dots, I_N . Тогда коэффициент роста погрешности v_A окажется конечным. Иными словами, при вычислении конечной серии интегралов алгоритм формально оказывается устойчивым. Тем не менее уже при не очень больших значениях N он настолько плохо обусловлен, что в практическом плане может считаться неустойчивым.

Для решения хорошо обусловленной задачи нет смысла применять плохо обусловленный алгоритм. Именно такими являются алгоритмы, первоначально предложенные в примерах 3.26 и 3.30.

Вернемся к примеру 3.30. Задача вычисления функции e^x хорошо обусловлена (см. пример 3.10). Можно ли было предвидеть катастрофическую потерю точности при вычислении значения e^{-8} прямым суммированием ряда (3.20)?

Рассмотрим задачу суммирования ряда $\sum_{k=0}^{\infty} a_k$ со слагаемыми $a_k = \frac{x^k}{k!}$. Каждое из этих слагаемых вычисляется с относительной погрешностью $\delta(a_k^*) \gtrsim \varepsilon_m$. При $x < 0$ формула (3.10) с учетом разложения (3.20) дает значение $v_\delta = \sum_{k=0}^{\infty} \frac{|x|^k}{k!} / \left| \sum_{k=0}^{\infty} \frac{x^k}{k!} \right| = e^{2|x|}$. Рост модуля x

¹ Иногда такие плохо обусловленные алгоритмы называют численно неустойчивыми.

(для $x < 0$) приводит к резкому ухудшению обусловленности вычислений. Для $x = -8.1$, как в примере 3.30, имеем $v_8 = e^{16.2} \approx 10^7$. Поэтому неудивительна полная потеря точности при вычислениях на 6-разрядном десятичном компьютере.

Рассмотрим теперь обусловленность алгоритма прямого вычисления по формуле (3.21). При малых x вычисленные приближения $c_1^* \approx \cos x$ и $c_2^* \approx \cos 2x$ содержат погрешности порядка ε_m . Поэтому $\bar{\Delta}(y^*) \sim 2\varepsilon_m$. Учитывая, что $y \approx 1.5x^2$, находим оценку границы относительной погрешности $\bar{\delta}(y^*) \sim x^{-2}\varepsilon_m$. Число обусловленности $v \sim x^{-2}$ быстро растет с уменьшением $|x|$. Уже при $|x| \lesssim \sqrt{\varepsilon_m}$ эта оценка дает $\delta(y^*) \gg 1$, т.е. говорит о катастрофической потере точности.

Если алгоритм, предназначенный для решения хорошо обусловленной задачи, оказался плохо обусловленным, то его следует признать неудовлетворительным и попытаться построить более качественный алгоритм. В примерах 3.30 и 3.31 это удалось сделать сравнительно легко.

Однако для плохо обусловленных задач дело обстоит иначе. Ключ к пониманию дает следующее высказывание [81]: «Если задача плохо обусловлена, то никакие усилия, потраченные на организацию изощренных вычислений, не могут дать правильных ответов, исключая случайности». Здесь требуется серьезное переосмысление постановки вычислительной задачи.

§ 3.6. Различные подходы к анализу ошибок

1. Прямой анализ ошибок. Общий эффект влияния ошибок обычно учитывают следующим образом. Изучают воздействие ошибок выходных данных, метода и округлений на получаемый результат y^* и пытаются оценить некоторую меру близости y^* к истинному решению y . Такой метод исследования называют *прямым анализом ошибок*. В большинстве случаев в данной книге мы будем следовать этому традиционному пути. Во многих (но далеко не во всех) случаях оценки погрешности удается получить; однако довольно часто они оказываются сильно завышенными и приводят к неоправданному пессимизму в оценке качества приближенного решения. Реальное значение погрешности $y - y^*$ часто значительно меньше, чем ее оценка, рассчитанная на самый неблагоприятный случай и выведенная с помощью прямого анализа. Особенно трудным является прямой анализ вычислительной погрешности.

2. Обратный анализ ошибок. В последнее время получил широкое распространение другой подход к оценке влияния ошибок. Оказывается, что довольно часто приближенное решение y^* можно трактовать как точное решение той же задачи, но отвечающее возмущенным входным данным x^* . Оценка такого эквивалентного возмущения входных данных и является целью *обратного анализа ошибок*.

В прикладных задачах входные данные, как правило, содержат погрешности. Обратный анализ показывает, что ошибки, внесенные в решение в процессе его вычисления, оказываются равносильными некоторым дополнительным ошибкам внесенным во входные данные. Сопоставление величины этого эквивалентного возмущения и уровня ошибок входных данных позволяет судить о качестве найденного решения. На рис. 3.5 представлена графическая иллюстрация обратного анализа ошибок. Здесь данное x^* таково, что решением задачи, соответствующим x^* , является y^* — результат приближенного решения задачи с входным данным x . На рисунке заштрихована *область неопределенности* входного данного; в пределах этой области входные данные для решающего задачу неразличимы. В представленном случае x^* оказалось внутри этой области, поэтому результат y^* следует признать вполне приемлемым.

Каждый из указанных двух подходов к оценке погрешности имеет свои достоинства и полезным является разумное их сочетание.

Пример 3.32. Пусть на 6-разрядном десятичном компьютера вычисляется корень уравнения

$$x^5 - 12.5x^4 + 62.5x^3 - 156.25x^2 + 195.3125x - 97.65625 = 0.$$

При вводе в компьютер последние два коэффициента будут округлены до шести значащих цифр и уравнение примет вид:

$$P^*(x) = x^5 - 12.5x^4 + 62.5x^3 - 156.25x^2 + 195.313x - 97.6563 = 0. \quad (3.22)$$

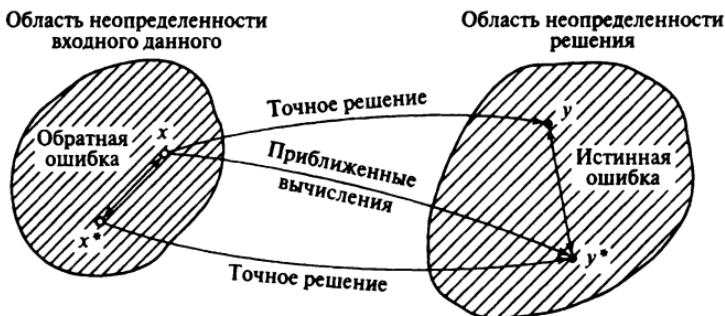


Рис. 3.5

Допустим, что некоторый алгоритм, примененный к этому уравнению, дал значение приближенного решения $x^* = 2.6$. Следуя логике прямого анализа ошибок, чтобы оценить качество полученного приближения нужно было бы задаться вопросом: насколько отличается x^* от истинного решения \bar{x} уравнения (3.22)? Мы поступим иначе. Подставляя x^* в левую часть уравнения (3.22), получаем значение $P^*(x^*) = 0.00126$. Заметим теперь, что $x^* = 2.6$ является точным решением уравнений

$$x^5 - 12.5x^4 + 62.5x^3 - 156.25x^2 + 195.3125x - 97.65626 = 0;$$

$$x^5 - 12.5x^4 + 62.5x^3 - 156.25x^2 + 195.31252x - 97.656312 = 0,$$

которые после округления коэффициентов до шести значащих цифр совпадают с уравнением (3.22) и становятся неотличимы от исходного уравнения. Найденное решение следует признать превосходным с точки зрения «философии» обратного анализа ошибок, так как оно является точным решением задачи, лежащей в пределах области неопределенности задачи (3.22). Рассчитывать на то, что после записи уравнения в виде (3.22) удастся получить лучший ответ, просто бессмысленно. Конечно, это немного обидно, особенно если учесть, что исходное уравнение в действительности есть развернутая запись уравнения $(x - 2.5)^5 = 0$ и истинным значением корня является $x = 2.5$. ▲

Представляется, что значение обратного анализа ошибок недостаточно осознанно, в особенности среди непрофессиональных вычислителей. Этот подход показывает на возможность иного взгляда на оценку качества приближенного решения, а значит, и на качество многих вычислительных алгоритмов. Сложившийся стереотип заставляет искать приближенное решение y^* математической задачи, мало отличающееся от ее истинного решения y . Однако для большинства практических задач в силу неопределенности в постановке и входных данных в действительности существует и область неопределенности решения (рис. 3.5). Поскольку эта область неизвестна, оценить степень близости вычисленного решения y^* к ней очень трудно. Гораздо проще, быть может, получить ответ на аналогичный вопрос для входного данного x^* , соответствующего решению y^* . Поэтому можно сформулировать цель вычислений и так: «найти точное решение задачи, которая мало отличается от поставленной задачи» или же так: «найти решение задачи с входным данным x^* , находящимся в пределах области неопределенности заданного входного данного x ».

Пример 3.33. Пусть для задачи Коши

$$y' - \sqrt{y} = \alpha(x)x, \quad y(1) = 2, \quad (3.23)$$

где $\alpha(x) = 2.6 + 0.01 e^{-x}$, найдено приближенное решение $y = 2x^2$. Как оценить его качество?

Условие $y(1) = 2$, очевидно, выполнено. Подставляя $y(x)$ в левую часть уравнения, убеждаемся, что $y(x)$ удовлетворяет уравнению (3.23) с коэффициентом $\alpha^* = 4 - \sqrt{2} \approx 2.586$ вместо α . Обратим внимание на то, что «истинное» значение α_0 коэффициента α нам в действительности неизвестно и $\alpha(x)$ — лишь некоторое его приближение. Числовой параметр 2.6 в лучшем случае получен округлением «истинного» значения, и, следовательно, α может отличаться от α_0 и на 0.05. Так как $|\alpha^* - \alpha| \leq 0.03$, то в силу естественной неопределенности в постановке задачи функция $y = 2x^2$ с позиции обратного анализа ошибок может считаться таким же равноправным решением поставленной задачи, как и найденное сколь угодно точно решение задачи (3.23). Во всяком случае теперь для того, чтобы отказаться от найденного приближенного решения, нужны довольно веские аргументы. ▲

Подчеркнем, что в основе методов решения некорректных и плохо обусловленных задач также лежит существенное переосмысление постановок вычислительных задач.

3. Статистический анализ ошибок. Даже для простых алгоритмов строгий анализ влияния ошибок округления очень сложен. При большом числе выполняемых операций гарантированные оценки погрешности, рассчитанные на самый неблагоприятный случай, как правило, бывают сильно завышенными.

Можно надеяться на то, что появляющиеся в реальном вычислительном процессе ошибки округления случайны и их взаимное влияние приводит к определенной компенсации результирующей ошибки. Статистический анализ ошибок исходящий из предположения об их случайности, направлен на исследование не максимально возможных, а наиболее вероятных ошибок.

Для сравнения покажем отличие в результатах на примере задачи вычисления суммы $S_N = \sum_{k=1}^N a_k$ большого числа положительных слагаемых. Гарантированная оценка погрешности дает значение относительной погрешности $\delta(S_N^*)$, растущее пропорционально N . В то же время статистический анализ показывает, что если ошибки округления являются случайными с нулевым средним значением, то $\delta(S_N^*)$ растет пропорционально \sqrt{N} , т.е. гораздо медленнее. К сожалению, на тех компь-

ютерах, где округление производится усечением, последнее предположение не выполнено, так как ошибки округления смешены в одну сторону и поэтому имеют ненулевое среднее значение. Здесь $\delta(S_N^*)$ растет опять пропорционально N .

4. Некоторые нестрогие способы анализа. Распространенным методом оценки влияния вычислительной погрешности является расчет с обычной и удвоенной точностью. Если результаты двух вычислений получаются существенно различными, это является свидетельством плохой обусловленности алгоритма. В то же время есть надежда на то, что совпадающие в ответах цифры верны.

Примерно такие же суждения о чувствительности решения к ошибкам округления можно сделать, если провести вычисления на двух различных вычислительных машинах или использовать различные компиляторы. Разумеется, такое исследование имеет смысл провести на одном-двух типичных примерах до начала массовых однотипных расчетов.

Влияние ошибок во входных данных на результат вычислений можно увидеть, если решить задачу несколько раз, изменив случайным образом входные данные в пределах ошибки их задания. Проделав такой эксперимент несколько раз, можно грубо оценить погрешность решения, вызванную погрешностями входных данных.

§ 3.7. Требования, предъявляемые к вычислительным алгоритмам

В § 3.4 и 3.5 были сформулированы два важнейших требования — корректность и хорошая обусловленность. Кроме них, к алгоритмам предъявляется еще целый ряд существенных требований.

1. Требования к абстрактным алгоритмам. К числу этих требований относятся: экономичность, надлежащая точность, экономия памяти, простота.

Экономичность алгоритма измеряется числом элементарных операций, необходимых для его реализации, и в конечном итоге сводится к затратам машинного времени. Это требование формулируют иногда как требование максимальной быстроты исполнения алгоритма. Экономичность особенно важна при массовых расчетах. Естественно, что при создании алгоритмов большое внимание уделяют минимизации числа операций. Для некоторых задач разработаны алгоритмы, требующие минимально возможного числа операций. Пример такого алго-

ритма — схема Горнера (см. пример 3.19). Отметим, что ряд математических алгоритмов, созданных в домашний период, оказался удивительно незакономичным. Приведем классический пример такого алгоритма.

Пример 3.34 (правило Крамера¹). Для решения системы линейных алгебраических уравнений $Ax = b$ порядка m по правилу Крамера предлагается вычислять компоненты вектора x как отношения специальным образом построенных определителей: $x_i = \Delta_i / \Delta$, $i = 1, 2, \dots, m$. Если вычислять определитель непосредственно по его определению, то нужно выполнить $(m - 1)m!$ умножений и $m!$ сложений. Пусть для вычислений используется компьютер с производительностью 10^6 умножений в секунду и решается система с числом неизвестных $m = 15$, весьма скромным для приложений. Тогда вычисление только одного определителя потребует $14 \cdot 15! \approx 1.8 \cdot 10^{13}$ умножений, в результате чего на вычисление решения уйдет около 10 лет непрерывной работы компьютера. Вместе с тем для решения той же системы на том же компьютере методом Гаусса (см. гл. 5) потребуется примерно 0.002 с. Естественно, что как вычислительный алгоритм правило Крамера следует забраковать. ▲

Даже для самых простых задач выбор экономичного алгоритма может дать существенное уменьшение числа операций.

Пример 3.35. Пусть требуется вычислить x^n , где n — натуральное число. Вычисление этой величины последовательным умножением на x предполагает выполнение $n - 1$ операций умножения. Нетрудно убедиться в том, что этот способ не самый экономичный. Например, x^{64} можно найти, выполнив не 63, а всего шесть операций умножения, если последовательным возведением в квадрат вычислить $x^2, x^4, x^8, x^{16}, x^{32}, x^{64}$.

В общем случае представим n в виде разложения (2.25) по степеням числа 2 (именно так число n хранится в памяти компьютера). Тогда

$$x^n = (x^{2^L})^{\alpha_L} \cdot (x^{2^{L-1}})^{\alpha_{L-1}} \cdot (x^2)^{\alpha_1} \cdot x^{\alpha_0}. \quad (3.24)$$

Заметим, что в произведении (3.24) следует учитывать при вычислении только те сомножители, для которых $\alpha_i \neq 0$ (т.е. $\alpha_i = 1$). Алгоритм,

¹ Габриэль Крамер (1704—1752) — швейцарский математик.

основанный на разложении (3.24), называется *бинарным алгоритмом*. Он позволяет найти x^n не более чем за $2 \log_2 n$ операций умножения. ▲

Замечание. Традиционно временные затраты на выполнение алгоритма оцениваются числом выполняемых им операций с плавающей точкой (флопов). Нужно иметь в виду, что на современных компьютерах такие оценки не вполне корректны, поскольку игнорируют временные затраты на передачу данных. В первую очередь это относится к параллельным компьютерам, но верно также и для рабочих станций и персональных компьютеров. Временные затраты на пересылку данных в персональных компьютерах могут превысить время на выполнение арифметических операций.

Требование *точности* означает, что вычислительный алгоритм должен давать решение задачи с заданной или приемлемой для задачи точностью ε .

Важным является требование *экономии памяти*. Хотя в последнее время доступная память компьютеров существенно расширилась, для «больших» задач требование экономии памяти может в ряде случаев стать основным. Интерес к экономическому размещению информации в памяти возрастает в связи со все более широким использованием персональных компьютеров для решения научно-технических и инженерных задач.

Учитывая необходимость дальнейшей программной реализации алгоритма, подчеркнем, что *простота алгоритма* также является весьма желательным его свойством.

2. Требования к программным реализациям алгоритмов. К настоящему времени выработан ряд требований к программам, реализующим вычислительные алгоритмы и предназначенным для длительного и широкого использования. Перечислим некоторые из них: надежность, работоспособность (робастность), переносимость (портабельность), поддерживаемость, простота в использовании и др. Рассмотрим эти требования более подробно.

Надежность программы означает, что она не содержит ошибок и вычисляет именно тот результат, для которого она предназначена.

Работоспособность (робастность) включает в себя надежность и предполагает, что программа способна выявлять недопустимые исходные данные, обнаруживать различные критические для задачи или алгоритма ситуации. Робастная программа реагирует на такие ситуации приемлемым для пользователя образом. Она составлена так, чтобы исключить недо-

пустимые операции типа деления на нуль и попытки извлечь квадратный корень или вычислить логарифм от отрицательного числа.

Алгоритм может «потерпеть неудачу» при решении задачи, если заданное входное данное не является для него допустимым. Конечно, в простых ситуациях пользователь должен сам различать допустимые для алгоритма входные данные от недопустимых. Однако чаще всего сделать это до вычислений очень трудно или невозможно, и в программе должен быть предусмотрен анализ данных и сообщение пользователю о недопустимых или сомнительных данных. Необходимо исключить ситуацию, характерную для некачественных программ, когда реакцией на задание данных, при которых алгоритм не может по объективным причинам найти решение задачи, является прекращение работы программы без сообщения причины останова, или же выдача внешне вполне правдоподобного, но совершенно бессмысленного результата.

Переносимость (портабельность) означает, что программа может работать на различных компьютерах без изменения или с незначительными изменениями. Всякая характеристика компьютера, используемая в программе (например, значение машинного эпсилон ϵ_m), должна или вычисляться самой программой, или задаваться пользователем.

Замечание. Принятие IEEE стандарта значительно упростило решение проблемы переносимости.

Поддерживаемость означает прежде всего требование легкости модификации. Для того чтобы была возможность внесения в программу изменений с минимальной вероятностью появления ошибок, она должна быть составлена максимально ясно и логично. Полезно вносить в текст программы содержательные комментарии. Разобраться в плохо составленной программе может оказаться труднее, чем создать новую. Поддерживаемая программа должна быть хорошо документирована. Плохое описание программы в лучшем случае способно вызвать к ней недоверие, а в худшем — может не позволить пользователю правильно ее эксплуатировать. К сожалению, нередка ситуация, когда предназначеннная для широкого использования программа настолько плохо документирована, что пользователь предпочитает потратить время на написание аналогичной программы (возможно, гораздо худшего качества) либо вообще отказаться от решения задачи.

Простота в использовании программы — весьма желательное, но трудно достижимое свойство. Часто добиться простоты в использовании можно только жертвуя надежностью или экономичностью. Существует ряд широко используемых программ, которые, в первую очередь, популярны, благодаря простоте в использовании.

3. Противоречивость требований. Можно продолжить перечисление требований к вычислительным алгоритмам, добавив, например требования универсальности и гибкости. Однако нетрудно понять, что сформулированные требования противоречивы. Большинство из них вступает в противоречие с экономичностью, выраженной через затраты машинного времени. В разных ситуациях на первый план может выступать то или иное требование и, удовлетворяя в большей степени одним требованиям, программа с неизбежностью в меньшей степени удовлетворяет другим. Это частично объясняет наличие большого числа программ, предназначенных для решения одной и той же задачи.

Естественно, что хорошая программа, которую можно предъявить для широкого использования, не может быть простой. Следует признать, что составление таких программ — это работа, требующая высокой квалификации и специальных знаний. Ее выполняют специалисты по созданию математического обеспечения компьютеров. Рядовой пользователь должен по возможности стремиться максимально использовать стандартные программы, а не создавать новые.

§ 3.8. Дополнительные замечания

1 Учебные пособия [97, 37] можно рассматривать как введение в теорию методов решения некорректных задач. Их удачно дополняют следующие книги [7, 72, 98]. Отметим, что [98] содержит не только теорию и алгоритмы, но и тексты соответствующих программ.

2 Весьма содержательное изложение проблемы обусловленности вычислительных задач и обратного анализа ошибок содержится в [81]. Здесь же подробно обсуждаются проблема создания высококачественного математического обеспечения и требования, предъявляемые к вычислительным алгоритмам.

Глава 4

МЕТОДЫ ОТЫСКАНИЯ РЕШЕНИЙ НЕЛИНЕЙНЫХ УРАВНЕНИЙ

В этой главе рассматривается задача отыскания корней нелинейных уравнений и излагаются методы ее решения. Это делается несколько подробнее, чем обычно принято в учебниках по численным методам. Дело в том, что нелинейное уравнение представляет собой редкий пример задачи, которая может быть сравнительно полно исследована элементарными средствами и допускает наглядные геометрические иллюстрации. В то же время многие проблемы, возникающие при отыскании корней нелинейных уравнений, типичны, а некоторые методы их решения (в особенности метод простой итерации и метод Ньютона) допускают широкие обобщения и играют в вычислительной математике фундаментальную роль.

§ 4.1. Постановка задачи. Основные этапы решения

1. Постановка задачи. Задача отыскания корней нелинейного уравнения с одним неизвестным вида

$$f(x) = 0 \quad (4.1)$$

имеет многовековую историю, но не потеряла свою актуальность и в наши дни. Она часто возникает как элементарный шаг при решении различных научных и технических проблем. Напомним, что *корнем* (или *решением*) уравнения (4.1) называется значение \bar{x} , при котором $f(\bar{x}) = 0$.

Для справедливости большинства рассуждений данной главы достаточно предположить, что в окрестности каждого из искомых корней функция $f(x)$ дважды непрерывно дифференцируема.

Корень \bar{x} уравнения (4.1) называется *простым*, если $f'(\bar{x}) \neq 0$. В противном случае (т.е. в случае $f'(\bar{x}) = 0$) корень \bar{x} называется *кратным*. Натуральное число m назовем *кратностью корня* \bar{x} , если $f^{(k)}(\bar{x}) = 0$ для $k = 1, 2, \dots, m - 1$ и $f^{(m)}(\bar{x}) \neq 0$. Геометрически корень \bar{x} соответствует точке пересечения графика функции $y = f(x)$ с осью Ox . Корень \bar{x} является простым, если график пересекает ось Ox под ненулевым углом, и кратным, если пересечение происходит под нулевым

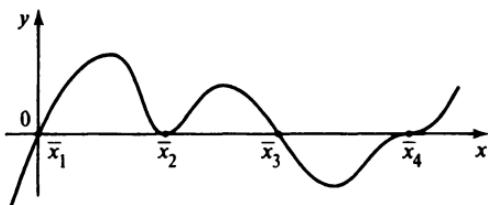


Рис. 4.1

углом. Функция $f(x)$, график которой изображен на рис. 4.1, имеет четыре корня. Корни \bar{x}_1 и \bar{x}_3 — простые, \bar{x}_2 и \bar{x}_4 — кратные.

Задача отыскания простых корней является существенно более простой (и чаще встречающейся), чем задача отыскания кратных корней. В действительности большинство методов решения уравнения (4.1) ориентировано именно на вычисление простых корней.

2. Уточнение постановки задачи. В конкретной задаче часто интерес представляют не все корни уравнения, а лишь некоторые из них. Тогда постановку задачи уточняют, указывая на то, какие из корней подлежат определению (положительные корни, корни из заданного интервала, максимальный из корней и т.д.).

В подавляющем большинстве случаев представить решение уравнения (4.1) в виде конечной формулы оказывается невозможным. Даже для простейшего алгебраического уравнения n -й степени

$$x^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0 = 0 \quad (4.2)$$

явные формулы, выражающие его корни через коэффициенты с помощью конечного числа арифметических операций и операций извлечения корней степени не выше n , найдены¹ лишь при $n = 2, 3, 4$. Однако уже для уравнений пятой и более высоких степеней таких формул не существует. Этот замечательный факт, известный как теорема Абеля, был установлен в 30-е годы XIX в. Н. Абелем² и Э. Галуа³.

¹ Алгебраические уравнения третьей и четвертой степеней не поддавались усилиям математиков около 2000 лет. Этую задачу решили итальянские математики эпохи Ренессанса Сципион дель Ферро (1456—1526), Никколо Тарталья (1500—1557), Джироламо Кардано (1501—1576), Людовико Феррари (1522—1565).

² Нильс Хенrik Абель (1802—1829) — норвежский математик. Теорема, о которой идет речь, была доказана им в возрасте около 22 лет.

³ Эварист Галуа (1811—1832) — французский математик, один из создателей теории групп. Убит на дуэли в возрасте 21 года. Краткий очерк о его жизни см. в [54].

Невозможность найти точное решение нелинейного уравнения кажется огорчительной. Однако нужно признать, что желание найти точное числовое значение решения вряд ли следует считать разумным. Во-первых, в реальных исследованиях зависимость $y = f(x)$ является лишь приближенным описанием, моделирующим истинную связь между параметрами y и x . Поэтому точное решение \bar{x} уравнения (4.1) все равно является лишь приближенным значением того параметра x , который в действительности соответствует значению $y = 0$. Во-вторых, даже если уравнение (4.1) допускает возможность нахождения решения в виде конечной формулы, то результат вычислений по этой формуле почти с неизбежностью содержит вычислительную погрешность и поэтому является приближенным.

Пример 4.1. Предположим, что исследование некоторого явления привело к необходимости решить уравнение

$$x^2 - 3.3x + 2.7 = 0. \quad (4.3)$$

Воспользовавшись формулами (3.2) для корней квадратного уравнения, получим значения $\bar{x}_1 = 1.5$, $\bar{x}_2 = 1.8$. Найдены ли нами точные значения параметра x ? Очевидно, нет. Скорее всего коэффициенты уравнения (4.3) известны приближенно и в лучшем случае они представляют округленные значения «истинных» коэффициентов. В действительности можно лишь утверждать, что $\bar{x}_1 \approx 1.5$, $\bar{x}_2 \approx 1.8$.

Предположим теперь, что «истинный» вид уравнения (4.3) таков: $x^2 - 3.3287x + 2.6631 = 0$. Тогда точные значения параметра можно вычислить по формуле $x_{1,2} = (3.3287 \pm \sqrt{3.3287^2 - 4 \cdot 2.6631})/2$. Однако она лишь указывает на то, какие операции и в каком порядке следует выполнить. В данном случае точное вычисление по формуле невозможно, так как она содержит операцию извлечения квадратного корня. Вычисленные по ней значения \bar{x}_1 , \bar{x}_2 неизбежно окажутся приближенными. ▲

В дальнейшем мы откажемся от попыток найти точные значения корней уравнения (4.1) и сосредоточим внимание на методах решения более реалистичной задачи приближенного вычисления корней с заданной точностью ε .

3. Основные этапы решения. Решение задачи отыскания корней нелинейного уравнения осуществляют в два этапа. Первый этап называется *этапом локализации* (или *отделения*) корней, второй — *этапом итерационного уточнения* корней.

Локализация корней. Отрезок $[a, b]$, содержащий только один корень \bar{x} уравнения (4.1), называют *отрезком локализации корня \bar{x}* . Цель этапа локализации считают достигнутой, если для каждого из подлежащих определению корней удалось указать отрезок локализации (его длину стараются по возможности сделать минимальной).

Прежде чем переходить непосредственно к отысканию отрезков локализации, имеет смысл провести предварительное исследование задачи для выяснения того, существуют ли вообще корни уравнения (4.1), сколько их и как они расположены на числовой оси.

Способы локализации корней многообразны, и указать универсальный метод не представляется возможным. Иногда отрезок локализации известен либо он определяется из физических соображений. В простых ситуациях хороший результат может давать графический метод (см. пример 4.2). Широко применяют построение таблиц значений функций f вида $y_i = f(x_i)$, $i = 1, 2, \dots, n$. При этом способе локализации о наличии на отрезке $[x_{i-1}, x_i]$ корня судят по перемене знака функции на концах отрезка (см. пример 4.3). Основанием для применения указанного способа служит следующая хорошо известная теорема математического анализа.

Теорема 4.1. Пусть функция f непрерывна на отрезке $[a, b]$ и принимает на его концах значения разных знаков, т.е. $f(a) \cdot f(b) < 0$. Тогда отрезок $[a, b]$ содержит по крайней мере один корень уравнения $f(x) = 0$. ■

К сожалению, корень четной кратности не удается локализовать на основании перемены знака с помощью даже очень подробной таблицы. Дело в том, что в малой окрестности такого корня (например, корня \bar{x}_2 на рис. 4.1) функция f имеет постоянный знак.

Важно подчеркнуть, что далеко не всегда для успешного отыскания корня \bar{x} уравнения (4.1) необходимо полное решение задачи локализации. Часто вместо отрезка локализации достаточно найти хорошее начальное приближение $x^{(0)}$ к корню \bar{x} .

Пример 4.2. Локализуем корни уравнения

$$4(1 - x^2) - e^x = 0. \quad (4.4)$$

Для этого преобразуем уравнение к виду $1 - x^2 = 0.25 e^x$ и построим графики функций $y = 1 - x^2$ и $y = 0.25 e^x$ (рис. 4.2). Абсциссы точек пересечения этих графиков являются корнями данного уравнения. Из рис. 4.2 видно, что уравнение имеет два корня \bar{x}_1 и \bar{x}_2 , расположенные на отрезках $[-1, 0]$ и $[0, 1]$. Убедимся, что функция $f(x) = 4(1 - x^2) - e^x$

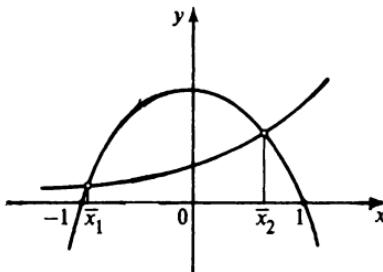


Рис. 4.2

принимает на концах указанных отрезков значения разных знаков. Действительно, $f(-1) = -e^{-1} < 0$, $f(0) = 3 > 0$, $f(1) = -e < 0$.

Следовательно, в силу теоремы 4.1 на каждом из отрезков $[-1, 0]$ и $[0, 1]$ находится по крайней мере один корень. ▲

Пример 4.3. Локализуем корни уравнения

$$x^3 - 1.1x^2 - 2.2x + 1.8 = 0.$$

Для этого составим таблицу значений функции $f(x) = x^3 - 1.1x^2 - 2.2x + 1.8$ на отрезке $[-2, 2]$ с шагом 0.4 (табл. 4.1).

Таблица 4.1						
x	-2.0	-1.6	-1.2	-0.8	-0.4	0.0
$f(x)$	-6.200	-1.592	1.128	2.344	2.440	1.800

Продолжение табл. 4.1						
x	0.4	0.8	1.2	1.6	2.0	
$f(x)$	0.808	-0.152	-0.696	-0.440	1.000	

Из таблицы видно, что функция f меняет знак на концах отрезков $[-1.6, -1.2]$, $[0.4, 0.8]$, $[1.6, 2.0]$. Теорема 4.1 дает основание утверждать, что каждый из этих отрезков содержит по крайней мере один корень. Учитывая, что в силу основной теоремы алгебры многочлен третьей степени не может иметь более трех корней, заключаем, что полученные три отрезка содержат ровно по одному корню. Таким образом, корни локализованы. ▲

Итерационное уточнение корней. На этом этапе для вычисления каждого из корней с точностью $\varepsilon > 0$ используют тот или иной итерационный метод, позволяющий построить последовательность $x^{(0)}, x^{(1)}, \dots, x^{(n)}, \dots$ приближений к корню \bar{x} .

Общее представление об итерационных методах и основные определения были даны в § 3.3. Введем дополнительно некоторые определения.

Итерационный метод называют *одношаговым*, если для вычисления очередного приближения $x^{(n+1)}$ используется только одно предыдущее приближение $x^{(n)}$ и *k-шаговым*, если для вычисления $x^{(n+1)}$ используются k предыдущих приближений $x^{(n-k+1)}, x^{(n-k+2)}, \dots, x^{(n)}$. Заметим, что для построения итерационной последовательности одношаговым методом требуется задание только одного начального приближения $x^{(0)}$, в то время как при использовании *k*-шагового метода — k начальных приближений $x^{(0)}, x^{(1)}, \dots, x^{(k-1)}$.

Скорость сходимости — одна из важнейших характеристик итерационных методов. Говорят, что *метод сходится со скоростью геометрической прогрессии*, знаменатель которой $0 < q < 1$, если для всех n справедлива следующая оценка:

$$|x^{(n)} - \bar{x}| \leq c_0 q^n. \quad (4.5)$$

Как нетрудно видеть, из оценки (4.5) действительно вытекает сходимость метода.

Пусть одношаговый итерационный метод обладает следующим свойством: существует σ -окрестность корня \bar{x} такая, что если приближение $x^{(n)}$ принадлежит этой окрестности, то справедлива оценка

$$|x^{(n+1)} - \bar{x}| \leq C|x^{(n)} - \bar{x}|^p, \quad (4.6)$$

где $C > 0$ и $p \geq 1$ — постоянные. В этом случае число p называют *порядком сходимости метода*. Если $p = 1$ и $C < 1$, то говорят, что метод обладает *линейной скоростью сходимости* в указанной σ -окрестности корня. Если $p > 1$, то принято говорить о *сверхлинейной скорости сходимости*. При $p = 2$ скорость сходимости называют *квадратичной*, а при $p = 3$ — *кубической*. При наличии оценки (4.6) у *k*-шагового метода (при $k > 1$) число p также будем называть *порядком сходимости метода*.

Лемма 4.1. Пусть одношаговый итерационный метод обладает линейной скоростью сходимости в некоторой σ -окрестности корня \bar{x} . Тогда при любом выборе начального приближения $x^{(0)}$ из σ -окрестности корня \bar{x} итерационная последовательность $x^{(n)}$ не выходит за пределы этой окрестности, метод сходится со скоростью геометрической прогрессии со знаменателем $q = C < 1$ и имеет место следующая оценка погрешности:

$$|x^{(n)} - \bar{x}| \leq q^n |x^{(0)} - \bar{x}|, \quad n \geq 0. \quad (4.7)$$

Доказательство. Заметим, что принадлежность $x^{(n)}$ окрестности $(\bar{x} - \sigma, \bar{x} + \sigma)$ является следствием неравенства (4.7). В самом деле,

так как $q < 1$, то $|x^{(n)} - \bar{x}| \leq |x^{(0)} - \bar{x}| < \sigma$. Сходимость $x^{(n)}$ к \bar{x} также вытекает из (4.7).

Справедливость неравенства (4.7) установим методом индукции.

При $n = 0$ оно переходит в очевидное: $|x^{(0)} - \bar{x}| \leq |x^{(0)} - \bar{x}|$.

Пусть неравенство (4.7) выполнено при $n = m - 1$. Тогда

$$|x^{(m)} - \bar{x}| \leq q|x^{(m-1)} - \bar{x}| \leq q^m|x^{(0)} - \bar{x}|,$$

т.е. неравенство выполнено и при $n = m$. ■

Лемма 4.2. Пусть одншаговый итерационный метод в некоторой σ -окрестности корня \bar{x} имеет p -й порядок сходимости, где $p > 1$. Пусть $\delta > 0$ таково, что $\delta \leq \sigma$ и $C\delta^{p-1} \leq 1$, где C — постоянная из неравенства (4.6). Тогда при любом выборе начального приближения $x^{(0)}$ из δ -окрестности корня \bar{x} итерационная последовательность $x^{(n)}$ не выходит за пределы этой окрестности, метод сходится и справедлива оценка

$$|x^{(n)} - \bar{x}| \leq C_1 q^{p^n}, \quad n \geq 0, \quad (4.8)$$

где $q = C_1^{-1} |x^{(0)} - \bar{x}|$, $C_1 = C^{-1/(p-1)}$.

Доказательство. Заметим, что принадлежность $x^{(n)}$ окрестности $(\bar{x} - \delta, \bar{x} + \delta)$ является следствием неравенства (4.8). В самом деле, так как $q < C_1^{-1}\delta \leq 1$ и $1 \leq p^n$, то из (4.8) вытекает, что $|x^{(n)} - \bar{x}| \leq C_1 q = |x^{(0)} - \bar{x}| < \delta$. Сходимость $x^{(n)}$ к \bar{x} также следует из (4.8).

Справедливость неравенства (4.8) установим методом индукции. При $n = 0$ оно переходит в очевидное: $|x^{(0)} - \bar{x}| \leq |x^{(0)} - \bar{x}|$. Пусть неравенство (4.8) выполнено при $n = m - 1$. Докажем, что оно верно и при $n = m$. Используя условие (4.6), получаем

$$|x^{(m)} - \bar{x}| \leq C|x^{(m-1)} - \bar{x}|^p \leq C(C_1 q^{p^{m-1}})^p = C_1 q^{p^m}. \blacksquare$$

С помощью доказанных лемм исследование сходимости итерационных методов сводится только к получению оценки (4.6).

§ 4.2. Обусловленность задачи вычисления корня

Пусть \bar{x} — корень уравнения (4.1), подлежащий определению. Будем считать, что входными данными для задачи вычисления корня \bar{x} являются значения $f(x)$ функции f в малой окрестности корня. Так как значения $f(x)$ будут вычисляться на компьютере по некоторой программе, то в действительности задаваемые значения являются прибли-

женными и их следует обозначать через $f^*(x)$. Погрешности в значениях $f^*(x)$ могут быть связаны не только с неизбежными ошибками округления, но и с использованием для вычисления значений функции f приближенных методов. К сожалению, нельзя ожидать, что в окрестности корня относительная погрешность $\delta(f^*)$ окажется малой. Достаточно обратиться к примерам 2.11 и 3.11, чтобы убедиться в том, что даже для чрезвычайно простых функций $y = 1 - x$ и $y = \sin x$ в окрестности корней $\bar{x} = 1$ и $\bar{x} = \pi$ значения этих функций не могут быть найдены с малой относительной погрешностью. Реально рассчитывать можно лишь на то, что малой окажется абсолютная погрешность вычисления значений функции.

Будем предполагать, что в достаточно малой окрестности корня выполняется неравенство $|f(x) - f^*(x)| < \bar{\Delta}$, где $\bar{\Delta} = \bar{\Delta}(f^*)$ — граница абсолютной погрешности. Сама погрешность корня ведет себя крайне нерегулярно и в первом приближении может восприниматься пользователем как некоторая случайная величина. На рис. 4.3, а представлена идеальная ситуация, отвечающая исходной математической постановке задачи, а на рис. 4.3, б — реальная ситуация, соответствующая вычислениям значений функции f на компьютере.

Если функция f непрерывна, то найдется такая малая окрестность $(\bar{x} - \bar{\varepsilon}, \bar{x} + \bar{\varepsilon})$ корня \bar{x} , имеющая радиус $\bar{\varepsilon} > 0$, в которой выполняется неравенство

$$|f(x)| \lesssim \bar{\Delta}. \quad (4.9)$$

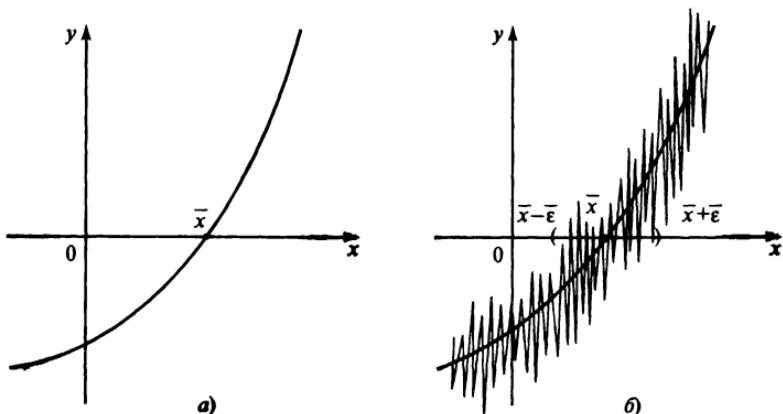


Рис.4.9

Для $x \in (\bar{x} - \bar{\varepsilon}, \bar{x} + \bar{\varepsilon})$ знак вычисленного значения $f^*(x)$, вообще говоря, не обязан совпадать со знаком $f(x)$ и, следовательно, становится невозможным определить, какое именно значение x из интервала $(\bar{x} - \bar{\varepsilon}, \bar{x} + \bar{\varepsilon})$ обращает функцию f в нуль (рис. 4.3, б).

Будем называть этот интервал *интервалом неопределенности корня \bar{x}* . Найдем оценку величины $\bar{\varepsilon}$. Пусть корень \bar{x} — простой. Для близких к \bar{x} значений x справедливо приближенное равенство

$$f(x) \approx f(\bar{x}) + f'(\bar{x})(x - \bar{x}) = f'(\bar{x})(x - \bar{x}).$$

Поэтому неравенство (4.9) имеет вид $|f'(\bar{x})(x - \bar{x})| \leq \bar{\Delta}$, откуда получаем

$$\bar{x} - \frac{\bar{\Delta}}{|f'(\bar{x})|} \lesssim x \lesssim \bar{x} + \frac{\bar{\Delta}}{|f'(\bar{x})|}.$$

Следовательно,

$$\bar{\varepsilon} \approx v_{\Delta} \bar{\Delta}(f^*). \quad (4.10)$$

Здесь $v_{\Delta} = \frac{1}{|f'(\bar{x})|}$ — число, которое в рассматриваемой задаче играет роль абсолютного числа обусловленности. Действительно, если \bar{x}^* — корень уравнения $f^*(x) = 0$, то \bar{x}^* удовлетворяет неравенству (4.9) и тогда выполнено неравенство

$$|\bar{x} - \bar{x}^*| \leq \bar{\Delta}(\bar{x}^*) \leq \bar{\varepsilon} \approx v_{\Delta} \bar{\Delta}(f^*). \quad (4.11)$$

Заметим, что радиус интервала неопределенности прямо пропорционален погрешности $\bar{\Delta}$ вычисления значения f . Кроме того, $\bar{\varepsilon}$ возрастает (обусловленность задачи ухудшается) с уменьшением $|f'(\bar{x})|$, т.е. с уменьшением модуля тангенса угла наклона, под которым график функции пересекает ось Ox (рис. 4.4, а, б).

Если же $f'(\bar{x}) = 0$ (т.е. корень \bar{x} — кратный), то формула (4.10) уже не верна. Пусть кратность корня равна m . Тогда в силу формулы Тейлора¹ справедливо приближенное равенство

$$f(x) \approx f(\bar{x}) + f'(\bar{x})(x - \bar{x}) + \frac{f''(\bar{x})}{2!} (x - \bar{x})^2 + \dots + \frac{f^{(m)}(\bar{x})}{m!} (x - \bar{x})^m,$$

¹ Брук Тейлор (1685—1731) — английский математик и философ. Широко известная формула разложения функции в степенной ряд была получена им в 1712 г.

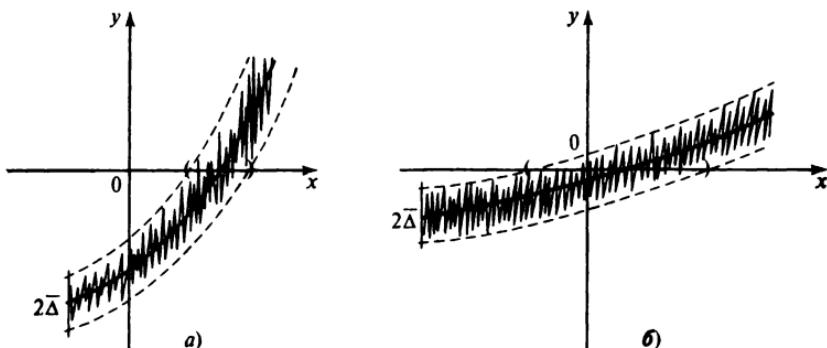


Рис. 4.4

в правой части которого все слагаемые, кроме последнего, равны нулю. Следовательно, неравенство (4.9) имеет вид

$$\left| \frac{f^{(m)}(\bar{x})}{m!} (x - \bar{x})^m \right| \lesssim \bar{\Delta} .$$

Решая его, получаем аналогично (4.10) оценку радиуса интервала неопределенности:

$$\bar{\epsilon} \approx \left| \frac{m!}{f^{(m)}(\bar{x})} \right|^{1/m} \bar{\Delta}^{1/m} .$$

Эта оценка означает, что для корня кратности m радиус интервала неопределенности пропорционален $\bar{\Delta}^{1/m}$, что свидетельствует о плохой обусловленности задачи вычисления кратных корней. С этой неприятностью мы уже сталкивались при рассмотрении примера 3.9.

Отметим, что $\bar{\epsilon}$ не может быть меньше $|\bar{x}| \epsilon_m$ — погрешности представления корня \bar{x} в компьютере.

В реальной ситуации оценить значение и даже порядок радиуса интервала неопределенности довольно сложно. Однако знать о его существовании нужно по крайней мере по двум причинам. Во-первых, не имеет смысла ставить задачу о вычислении корня \bar{x} с точностью $\epsilon < \bar{\epsilon}$. В условиях неопределенности, вызванных приближенным заданием функции, любое значение $\bar{x}^* \in (\bar{x} - \bar{\epsilon}, \bar{x} + \bar{\epsilon})$ может быть с одной и той же степенью достоверности принято за решение уравнения. Во-вторых, нельзя требовать от алгоритмов отыскания корня получения достоверных результатов после того, как очередное приближение попало в интервал неопределенности или оказалось очень близко от него, в этой

ситуации вычисления следует прекратить и считать, что получен максимум действительно возможного.

Для большинства итерационных методов определить этот момент можно, поскольку начиная с него поведение приближений $x^{(n)}$ становится крайне нерегулярным. Если вдали от интервала неопределенности величина

$$q^{(n)} = \frac{|x^{(n)} - x^{(n-1)}|}{|x^{(n-1)} - x^{(n-2)}|} \quad (4.12)$$

обычно бывает меньше единицы ($|x^{(n)} - x^{(n-1)}| < |x^{(n-1)} - x^{(n-2)}|$), то появление при некотором n значения $q^{(n)} > 1$ свидетельствует, скорее всего, о начале «разболтки» — хаотического поведения итерационной последовательности. В этой ситуации вычисления имеет смысл прервать, чтобы выяснить причину явления и принять правильное решение. Лучшим из полученных приближений к решению следует считать, конечно, $x^{(n-1)}$. Использование для контроля вычислений величины (4.12) называют часто *правилом Гарвика*.

§ 4.3. Метод бисекции

1. Описание метода. Пусть требуется с заданной точностью $\varepsilon > 0$ найти корень \bar{x} уравнения (4.1). Отрезок локализации $[a, b]$ (т.е. отрезок, содержащий только один корень \bar{x}) будем считать заданным. Предположим, что функция f непрерывна на отрезке $[a, b]$ и на его концах принимает значения разных знаков, т.е.

$$f(a)f(b) < 0. \quad (4.13)$$

На рис. 4.5 изображен случай, когда $f(a) < 0$ и $f(b) > 0$.

Для дальнейшего будет удобно обозначить отрезок $[a, b]$ через $[a^{(0)}, b^{(0)}]$. Примем за приближенное значение корня середину отрезка —

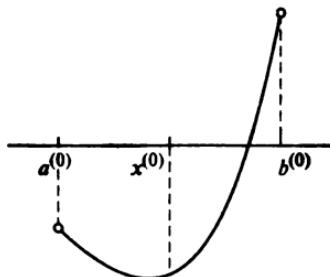


Рис. 4.5

точку $x^{(0)} = (a^{(0)} + b^{(0)})/2$. Так как положение корня \bar{x} на отрезке $[a^{(0)}, b^{(0)}]$ неизвестно, то можно лишь утверждать, что погрешность этого приближения не превышает половины длины отрезка (рис. 4.5):

$$|x^{(0)} - \bar{x}| \leq (b^{(0)} - a^{(0)})/2.$$

Уменьшить погрешность приближения можно, уточнив отрезок локализации, т.е. заменив начальный отрезок $[a^{(0)}, b^{(0)}]$ отрезком $[a^{(1)}, b^{(1)}]$ меньшей длины. Согласно методу бисекции (половинного деления) в качестве $[a^{(1)}, b^{(1)}]$ берут тот из отрезков $[a^{(0)}, x^{(0)}]$ и $[x^{(0)}, b^{(0)}]$, на концах которого выполняется условие $f(a^{(1)})f(b^{(1)}) \leq 0$. Этот отрезок содержит искомый корень. Действительно, если $f(a^{(1)})f(b^{(1)}) < 0$, то наличие корня следует из теоремы 4.1; если же $f(a^{(1)})f(b^{(1)}) = 0$, то корнем является один из концов отрезка. Середина полученного отрезка $x^{(1)} = (a^{(1)} + b^{(1)})/2$ дает теперь приближение к корню, оценка погрешности которого составляет

$$|x^{(1)} - \bar{x}| \leq (b^{(1)} - a^{(1)})/2 = (b - a)/2^2.$$

За очередное уточнение отрезка локализации $[a^{(2)}, b^{(2)}]$ снова берут тот из отрезков $[a^{(1)}, x^{(1)}]$, $[x^{(1)}, b^{(1)}]$, на концах которого выполняется условие $f(a^{(2)})f(b^{(2)}) \leq 0$.

Опишем очередную $(n+1)$ -ю итерацию метода. Пусть отрезок $[a^{(n)}, b^{(n)}]$ уже найден и вычислены значения $x^{(n)}$, $f(a^{(n)})$, $f(b^{(n)})$. Тогда производят следующие действия:

1°. Вычисляют $f(x^{(n)})$.

2°. Если $f(a^{(n)})f(x^{(n)}) \leq 0$, то в качестве отрезка локализации $[a^{(n+1)}, b^{(n+1)}]$ принимают отрезок $[a^{(n)}, x^{(n)}]$. В противном случае $f(x^{(n)})f(b^{(n)}) < 0$ и за $[a^{(n+1)}, b^{(n+1)}]$ принимают отрезок $[x^{(n)}, b^{(n)}]$.

3°. Вычисляют $x^{(n+1)} = (a^{(n+1)} + b^{(n+1)})/2$.

Неограниченное продолжение итерационного процесса дает последовательность отрезков $[a^{(0)}, b^{(0)}]$, $[a^{(1)}, b^{(1)}]$, ..., $[a^{(n)}, b^{(n)}]$, ..., содержащих искомый корень. Каждый из них (за исключением начального) получен делением пополам предыдущего отрезка.

2. Скорость сходимости. Середина n -го отрезка — точка $x^{(n)} = (a^{(n)} + b^{(n)})/2$ дает приближение к корню \bar{x} , имеющее оценку погрешности

$$|x^{(n)} - \bar{x}| \leq (b^{(n)} - a^{(n)})/2 = (b - a)/2^{n+1}. \quad (4.14)$$

Из этой оценки видно, что метод бисекции сходится со скоростью геометрической прогрессии, знаменатель которой $q = 1/2$. По сравнению с другими методами метод бисекции сходится довольно медленно.

Однако он очень прост и весьма непрятателен; для его применения достаточно, чтобы выполнялось неравенство (4.13), функция f была непрерывна и верно определялся ее знак. В тех ситуациях, где не нужна сверхвысокая скорость сходимости (а это часто имеет место при простых прикладных расчетах), этот метод весьма привлекателен.

Заметим, что число итераций, которое требуется при применении метода бисекции для достижения разумной точности ε , не может быть очень большим. Например, для уменьшения первоначального отрезка локализации в 10^6 раз нужно 19 итераций.

Замечание. Интересно, что при выполнении условия (4.13) метод бисекции сходится даже, если на отрезке $[a, b]$ находится не один, а несколько корней уравнения $f(x) = 0$. В результате применения метода будет найден один из этих корней.

3. Критерий окончания. Итерации следует вести до тех пор, пока не будет выполнено неравенство $b^{(n)} - a^{(n)} < 2\varepsilon$. При его выполнении в силу оценки (4.14) можно принять $x^{(n)}$ за приближение к корню с точностью ε .

Пример 4.4. Найдем методом бисекции с точностью $\varepsilon = 10^{-2}$ положительный корень уравнения $4(1 - x^2) - e^x = 0$.

В примере 4.2 этот корень был локализован на отрезке $[0, 1]$, причем $f(0) > 0, f(1) < 0$. Положим $a^{(0)} = 0, b^{(0)} = 1, x^{(0)} = (a^{(0)} + b^{(0)})/2 = 0.5$.

1-я итерация. Вычисляем $f(x^{(0)}) \approx 1.3512$. Так как $f(a^{(0)})f(x^{(0)}) > 0$, то за очередной отрезок локализации принимаем $[a^{(1)}, b^{(1)}] = [0.5, 1]$. Вычисляем $x^{(1)} = (a^{(1)} + b^{(1)})/2 = 0.75$.

2-я итерация. Вычисляем $f(x^{(1)}) \approx -0.3670$. Так как $f(a^{(1)})f(x^{(1)}) < 0$, то $[a^{(2)}, b^{(2)}] = [0.5, 0.75]$ и $x^{(2)} = (a^{(2)} + b^{(2)})/2 = 0.625$.

Результаты следующих итераций (с четырьмя цифрами после десятичной точки) приведены в табл. 4.2.

Таблица 4.2

Номер итерации k	$a^{(n)}$	$b^{(n)}$	Знак		$x^{(n)}$	$f(x^{(n)})$	$b^{(n)} - a^{(n)}$
			$f(a^{(n)})$	$f(b^{(n)})$			
0	0.0000	1.0000	+	-	0.5000	1.3513	1.0000
1	0.5000	1.0000	+	-	0.7500	-0.3670	0.5000
2	0.5000	0.7500	+	-	0.6250	0.5693	0.2500
3	0.6250	0.7500	+	-	0.6875	0.1206	0.1250
4	0.6875	0.7500	+	-	0.7187	-0.1182	0.0625
5	0.6875	0.7187	+	-	0.7031	0.0222	0.0312
6	0.7031	0.7187	+	-	0.7109	—	0.0156

При $n = 6$ имеем $b^{(6)} - a^{(6)} \approx 0.0156 < 2 \cdot 10^{-2}$. Следовательно, заданная точность достигнута и можно принять $\bar{x} \approx x^{(6)}$. Окончательно получим $\bar{x} = 0.71 \pm 0.01$. ▲

4. Влияние вычислительной погрешности. При использовании метода бисекции принципиально важным является правильное определение знака функции f . В случае, когда $x^{(n)}$ попадает в интервал неопределенности корня (см. § 4.2), знак вычисленного значения $f^*(x^{(n)})$ не обязан быть верным, и последующие итерации не имеют смысла. Однако этот метод следует признать очень надежным; он гарантирует точность приближения, примерно равную радиусу интервала неопределенности $\bar{\epsilon}$. Как было отмечено в § 4.3, большого требовать нельзя.

§ 4.4. Метод простой итерации

1. Описание метода. Чтобы применить метод простой итерации для решения нелинейного уравнения (4.1), необходимо преобразовать это уравнение к следующему виду:

$$x = \varphi(x). \quad (4.15)$$

Это преобразование (*приведение уравнения к виду, удобному для итераций*) можно выполнить различными способами; некоторые из них будут указаны ниже. Функцию φ далее будем называть *итерационной функцией*.

Выберем каким-либо образом приближенное значение корня $x^{(0)}$ и подставим его в правую часть уравнения (4.15). Получим значение $x^{(1)} = \varphi(x^{(0)})$. Подставляя теперь $x^{(1)}$ в правую часть уравнения (4.15), имеем $x^{(2)} = \varphi(x^{(1)})$. Продолжая этот процесс неограниченно, получаем последовательность приближений к корню, вычисляемых по формуле

$$x^{(n+1)} = \varphi(x^{(n)}), \quad n \geq 0. \quad (4.16)$$

Очевидно, что метод простой итерации — одношаговый (см. § 4.1).

Если существует предел построенной последовательности $\bar{x} = \lim_{n \rightarrow \infty} x^{(n)}$, то, переходя к пределу в равенстве (4.16) и предполагая функцию φ непрерывной, получаем равенство

$$\bar{x} = \varphi(\bar{x}). \quad (4.17)$$

Это значит, что \bar{x} — корень уравнения (4.15).

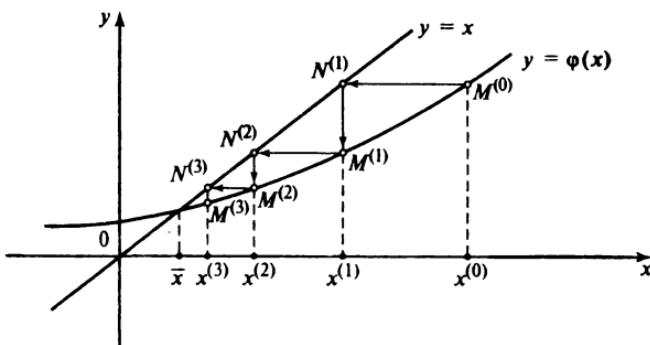


Рис. 4.6

2. Геометрическая иллюстрация. Из рис. 4.6 видно, что корень \bar{x} уравнения (4.15) является абсциссой точки пересечения графиков двух функций: $y = x$ и $y = \varphi(x)$. Возьмем некоторое начальное приближение $x^{(0)}$, которому отвечает расположенная на кривой $y = \varphi(x)$ точка $M^{(0)}$ с координатами $(x^{(0)}, x^{(0)})$ (напомним, что $x^{(1)} = \varphi(x^{(0)})$). Соединим точку $M^{(0)}$ отрезком прямой $y = x^{(1)}$ с лежащей на прямой $y = x$ точкой $N^{(1)}$ с координатами $(x^{(1)}, x^{(1)})$. Проведем теперь через точку $N^{(1)}$ прямую $x = x^{(1)}$ до пересечения с кривой $y = \varphi(x)$ в точке $M^{(1)}$ с координатами $(x^{(1)}, x^{(2)})$. Продолжая этот процесс далее, получаем ломаную линию $M^{(0)}N^{(1)}M^{(1)}N^{(2)}M^{(2)}\dots$, для которой абсциссы точек $M^{(n)}$ представляют собой последовательные приближения $x^{(n)}$ к решению \bar{x} .

3. Сходимость метода. На рис. 4.7, а—г представлена геометрическая иллюстрация поведения итерационного процесса в четырех простейших случаях взаимного расположения прямой $y = x$ и кривой $y = \varphi(x)$.

Рисунки 4.7, а и б иллюстрируют случаи, когда метод простой итерации сходится, причем, как нетрудно заметить, — при произвольном начальном приближении. Напротив, рисунки 4.7, в и г иллюстрируют случаи, когда метод расходится при любом выборе начального приближения $x^{(0)} \neq \bar{x}$. Заметим, что на рис. 4.7, а и б $|\varphi'(x)| < 1$ (модуль тангенса угла наклона кривой $y = \varphi(x)$ к оси абсцисс меньше единицы), а на рис. 4.7, в и г, наоборот, $|\varphi'(x)| > 1$. Таким образом, можно предположить, что сходимость метода простой итерации связана с выполнением условия $|\varphi'(x)| < 1$. Действительно, имеет место следующий результат.

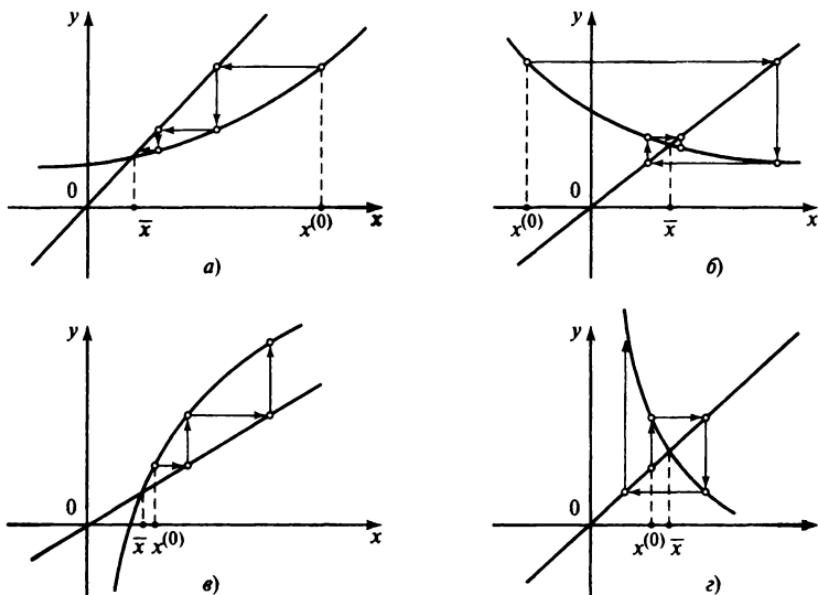


Рис. 4.7

Теорема 4.2. Пусть в некоторой σ -окрестности корня \bar{x} функция φ дифференцируема и удовлетворяет неравенству

$$|\varphi'(x)| \leq q, \quad (4.18)$$

где $0 \leq q < 1$ — постоянная.

Тогда независимо от выбора начального приближения $x^{(0)}$ из указанной σ -окрестности корня итерационная последовательность не выходит из этой окрестности, метод сходится со скоростью геометрической прогрессии и справедлива следующая оценка погрешности:

$$|x^{(n)} - \bar{x}| \leq q^n |x^{(0)} - \bar{x}|. \quad (4.19)$$

Доказательство. Вычитая из равенства (4.16) равенство (4.17) и используя формулу конечных приращений Лагранжа, получаем

$$x^{(n+1)} - \bar{x} = \varphi(x^{(n)}) - \varphi(\bar{x}) = \alpha^{(n+1)}(x^{(n)} - \bar{x}). \quad (4.20)$$

Здесь $\alpha^{(n+1)} = \varphi'(\xi^{(n)})$, где $\xi^{(n)}$ — некоторая точка, расположенная между $x^{(n)}$ и \bar{x} . Если $x^{(n)} \in (\bar{x} - \sigma, \bar{x} + \sigma)$, то $|\alpha^{(n+1)}| \leq q$ в силу условия (4.18). Тогда на основании равенства (4.20) получаем

$$|x^{(n+1)} - \bar{x}| \leq q |x^{(n)} - \bar{x}|.$$

Это означает, что метод простой итерации обладает линейной скоростью сходимости и поэтому доказательство теоремы завершается применением леммы 4.1. ■

Оценка погрешности (4.19) является априорной. Она показывает, что метод простой итерации сходится со скоростью геометрической прогрессии, знаменатель которой равен q . Чем меньше q , тем выше скорость сходимости. Видна и роль правильного выбора начального приближения: чем меньше погрешность начального приближения, тем меньше итераций потребуется сделать для достижения заданной точности ε .

Неравенство (4.19), как правило, не используется для практической оценки погрешности. Одна из причин этого состоит в том, что значение \bar{x} , входящее в правую часть оценки, неизвестно. Кроме того, использование неравенства (4.19) приводит к существенно завышенной оценке погрешности.

4. Критерий окончания. Выведем апостериорную оценку погрешности, пригодную для практического применения.

Теорема 4.3. Пусть выполнены условия теоремы 4.2 и $x^{(0)} \in (\bar{x} - \sigma, \bar{x} + \sigma)$. Тогда верна следующая апостериорная оценка погрешности:

$$|x^{(n)} - \bar{x}| \leq \frac{q}{1-q} |x^{(n)} - x^{(n-1)}|, \quad n \geq 1. \quad (4.21)$$

Доказательство. В силу равенства (4.20) имеем

$$x^{(n)} - \bar{x} = \alpha^{(n)}(x^{(n-1)} - \bar{x}) = \alpha^{(n)}(x^{(n-1)} - x^{(n)}) + \alpha^{(n)}(x^{(n)} - \bar{x}).$$

Откуда

$$x^{(n)} - \bar{x} = \frac{\alpha^{(n)}}{1-\alpha^{(n)}} (x^{(n-1)} - x^{(n)}). \quad (4.22)$$

Взяв модуль от левой и правой частей этого равенства и воспользовавшись неравенством $\left| \frac{\alpha^{(n)}}{1-\alpha^{(n)}} \right| \leq \frac{q}{1-q}$, получим требуемое соотношение (4.21). ■

Если значение q известно, то неравенство (4.21) дает эффективный метод контроля погрешности и можно сформулировать следующий критерий окончания итерационного процесса. Вычисления следует вести до выполнения неравенства $\frac{q}{1-q} |x^{(n)} - x^{(n-1)}| < \varepsilon$ или равносильного ему неравенства

$$|x^{(n)} - x^{(n-1)}| < \frac{1-q}{q} \varepsilon. \quad (4.23)$$

Если это условие выполнено, то найденное значение $x^{(n)}$ является приближением к \bar{x} с точностью ε .

Пример 4.5. Используем метод простой итерации для вычисления положительного корня \bar{x} уравнения $4(1 - x^2) - e^x = 0$ с точностью $\varepsilon = 10^{-4}$. Результат примера 4.4 дает для корня отрезок локализации $[a, b] = [0.70, 0.72]$.

Преобразуем уравнение к виду (4.15), где $\varphi(x) = \sqrt{1 - e^x}/4$. Заметим, что $\varphi'(x) = -e^x/(8\sqrt{1 - e^x}/4)$, $\varphi''(x) = -e^x(1 - e^x/8)/(8\sqrt{1 - e^x}/4)^3$. Так как $\varphi'' < 0$ на $[a, b]$, то производная φ' монотонно убывает и $q = \max_{[a, b]} |\varphi'(x)| = \varphi'(b) \approx 0.37$. Следовательно, условие сходимости (4.18) выполнено. Возьмем $x^{(0)} = 0.7$ и будем вести итерации до выполнения критерия (4.23). В табл. 4.3 соответствующие приближения приведены с десятью знакамиmantиссы.

Таблица 4.3

n	$x^{(n)}$	$\frac{q}{1-q} x^{(n)} - x^{(n-1)} $
0	0.7000000000	—
1	0.7046714292	$3 \cdot 10^{-3}$
2	0.7029968319	$1 \cdot 10^{-3}$
3	0.7035984939	$4 \cdot 10^{-4}$
4	0.7033824994	$2 \cdot 10^{-4}$
5	0.7034600632	$5 \cdot 10^{-5}$

Критерий окончания выполняется при $n = 5$. После округления значения $x^{(5)}$ до четырех значащих цифр получим $\bar{x} = 0.7035 \pm 0.0001$. ▲

Замечание. Часто в практике вычислений вместо критерия (4.23) используется привлекательный своей простотой критерий

$$|x^{(n)} - x^{(n-1)}| < \varepsilon. \quad (4.24)$$

В случае $0 < q \leq 1/2$ использование критерия (4.24) оправдано. Действительно, здесь $(1 - q)/q \geq 1$ и поэтому выполнение неравенства (4.24) влечет за собой выполнение неравенства (4.23).

В то же время при $1/2 < q < 1$ использование критерия (4.24) может привести к преждевременному прекращению итераций. Дело в том, что когда значение величины q близко к единице, итерационный процесс сходится медленно и расстояние между двумя последовательными приближениями $x^{(n)}$ и $x^{(n-1)}$ не характеризует расстояние от $x^{(n)}$ до решения \bar{x} (рис. 4.8).

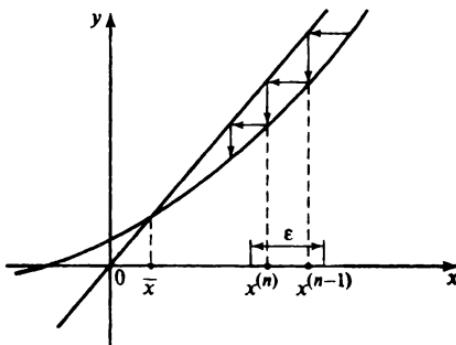


Рис. 4.8

Пример 4.6. Пусть метод простой итерации используется для решения уравнения $x = 0.9999x + 10^{-4}/\sqrt{2}$. Здесь $\phi(x) = 0.9999x + 10^{-4}/\sqrt{2}$, $\phi'(x) = 0.9999$ и, следовательно, условие (4.18) сходимости метода выполнено. Для вычислений по формуле (4.16) используем 6-разрядный десятичный компьютер. Возьмем $x^{(0)} = 0.715010$. Тогда $x^{(1)} = 0.715009$ и, если доверять критерию (4.24), то следовало бы считать, что решение получено с точностью $\epsilon = 10^{-6}$. Продолжая вычисления, имеем $x^{(2)} = 0.715008$, $x^{(3)} = 0.715007$, $x^{(4)} = 0.715006$, $x^{(5)} = 0.715005$, $x^{(6)} = 0.715005$. Дальнейшие итерации теряют смысл. Сколько же верных знаков найдено? Сравнение с точным значением решения $\bar{x} = 1/\sqrt{2} = 0.707106 \dots$ показывает, что верными в приближении $x^{(6)}$ являются только две значащие цифры. Использование критерия (4.24) в данном случае категорически недопустимо. В дальнейшем мы еще вернемся к обсуждению этого примера. ▲

Использование критерия (4.23) предполагает знание значения величины q , входящей в условие (4.18). Однако далеко не всегда это значение известно, либо может быть легко вычислено. В тех же случаях, когда удается оценить q , эта оценка оказывается довольно грубой.

Исключим из критерия окончания итераций величину q . Заметим, что в малой окрестности корня значение производной ϕ' практически постоянно: $\phi'(x) \approx \phi'(\bar{x})$. Поэтому в равенстве (4.22) величину $\alpha^{(n)} = \phi'(\xi^{(n-1)})$ можно приближенно заменить на $\phi'(\bar{x})$. Далее в силу равенства

$$x^{(n)} - x^{(n-1)} = \phi(x^{(n-1)}) - \phi(x^{(n-2)}) = \phi'(\tilde{\xi}^{(n)})(x^{(n-1)} - x^{(n-2)}),$$

где $\tilde{\xi}^{(n)}$ — промежуточная между $x^{(n-1)}$ и $x^{(n-2)}$ точка, имеем

$$\tilde{\alpha}^{(n)} = \frac{x^{(n)} - x^{(n-1)}}{x^{(n-1)} - x^{(n-2)}} = \varphi'(\tilde{\xi}^{(n)}) \approx \varphi'(\bar{x}).$$

Таким образом, в равенстве (4.22) можно положить $\alpha^{(n)} \approx \tilde{\alpha}^{(n)}$ и поэтому при определенных условиях можно использовать следующий практический критерий окончания итерационного процесса:

$$|x^{(n)} - x^{(n-1)}| \leq \left| \frac{1 - \tilde{\alpha}^{(n)}}{\tilde{\alpha}^{(n)}} \right| \varepsilon.$$

5. Дополнительные сведения о характере сходимости. Когда производная φ' знакопостоянна на отрезке локализации, итерационная последовательность обладает некоторыми полезными дополнительными свойствами. Заметим, что при $\varphi'(\bar{x}) \neq 0$ в достаточно малой окрестности корня знак производной $\varphi'(x)$ действительно постоянен.

Теорема 4.4. Пусть $[a, b]$ — отрезок локализации корня уравнения (4.15). Предположим, что на этом отрезке функция φ непрерывно дифференцируема, а ее производная $\varphi'(x)$ знакопостоянна и удовлетворяет неравенству (4.18) при $0 < q < 1$. Пусть $x^{(0)} \in [a, b]$ — произвольное начальное приближение и при $\varphi' < 0$ выполнено дополнительное условие $x^{(1)} = \varphi(x^{(0)}) \in [a, b]$ (первое приближение не выходит за пределы отрезка локализации).

Тогда итерационная последовательность не выходит за пределы отрезка $[a, b]$, метод простой итерации сходится и верны оценки погрешности (4.19), (4.21). Кроме того, справедливы следующие свойства:

1°. Если $\varphi' > 0$ на $[a, b]$, то $x^{(n)}$ сходится к \bar{x} , монотонно возрастая при $a \leq x^{(0)} < \bar{x}$, и монотонно убывая при $\bar{x} < x^{(0)} \leq b$.

2°. Если $\varphi' < 0$ на $[a, b]$, то сходимость $x^{(n)}$ к \bar{x} носит колебательный характер, т.е. при всех $n \geq 0$ значения $x^{(n)}$ и $x^{(n+1)}$ расположены по разные стороны от \bar{x} , причем последовательности приближений с четными и нечетными номерами сходятся к \bar{x} монотонно. В этом случае верна апостериорная оценка погрешности

$$|x^{(n)} - \bar{x}| \leq |x^{(n)} - x^{(n-1)}|, \quad n \geq 1 \quad (4.25)$$

и справедлив критерий (4.24) окончания итерационного процесса.

Мы не будем приводить здесь полное доказательство теоремы. Оно основано на использовании равенства (4.20), установленного при доказательстве теоремы 4.2. Докажем только справедливость свойств 1° и 2°.

Доказательство. Если $0 < \varphi' \leq q$, то знаки величин $x^{(n)} - \bar{x}$ и $x^{(n-1)} - \bar{x}$ совпадают, в то время как $|x^{(n)} - \bar{x}| \leq q|x^{(n-1)} - \bar{x}|$. Следовательно, последовательность монотонно приближается к \bar{x} с той стороны, где лежит $x^{(0)}$. Если же $-q < \varphi' < 0$, то из равенства (4.20) следует, что знаки величин $x^{(n)} - \bar{x}$ и $x^{(n-1)} - \bar{x}$ различны. Это подтверждает колебательный характер сходимости. ■

Замечание. Монотонный и колебательный характер сходимости итерационной последовательности, указанные в теореме 4.4, иллюстрируют соответственно рис. 4.7, а и б.

6. Приведение уравнения к виду, удобному для итераций. Ключевой момент в применении метода простой итерации — эквивалентное преобразование уравнения (4.1) к виду (4.15). Конечно, такое преобразование имеет смысл только тогда, когда оказывается выполненным условие (4.18) с $0 < q < 1$. Укажем один из простых способов такого преобразования.

Предположим, что производная f' на отрезке $[a, b]$ непрерывна и положительна. Тогда существуют положительные постоянные m и M такие, что $0 < m \leq f'(x) \leq M$, $x \in [a, b]$. Приведем уравнение (4.1) к виду

$$x = x - \alpha f(x), \quad (4.26)$$

где $\alpha > 0$. Здесь итерационная функция φ имеет вид $\varphi(x) = x - \alpha f(x)$. Как выбрать α , чтобы выполнялось условие (4.18), причем q было бы по возможности минимальным?

Заметим, что $1 - \alpha M \leq \varphi'(x) = 1 - \alpha f'(x) \leq 1 - \alpha m$ и поэтому $|\varphi'(x)| \leq q(\alpha) = \max\{|1 - \alpha M|, |1 - \alpha m|\}$. Для того чтобы было выполнено неравенство $q(\alpha) < 1$, достаточно взять любое $\alpha \in (0, 2/M)$. Конкретный выбор параметра α зависит от наличия информации о числах m и M . Если известны обе эти величины, то лучшим является выбор $\alpha = \alpha_0 = 2/(M+m)$. В этом случае $q(\alpha_0) = (M-m)/(M+m)$. Если же известно только M , то можно положить $\alpha = \alpha_1 = 1/M$. В этом случае $q(\alpha_1) = 1 - m/M$. Кроме того, при $\alpha = \alpha_1$ производная $\varphi'(x)$ неотрицательна, и в силу теоремы 4.3 сходимость является монотонной.

Замечание. Случай, когда производная f' отрицательна, сводится к рассмотренному выше умножением уравнения $f(x) = 0$ на -1 .

Пример 4.7. Для решения задачи, поставленной в примере 4.5, можно воспользоваться преобразованием уравнения (4.4) к виду (4.26). Будем считать известным, что $x \in [0.70, 0.72]$. Так как на отрезке локализации выполнено условие $(4(1-x^2) - e^x)' = -8x - e^x < 0$,

то перепишем уравнение в виде $f_1(x) = 0$, где $f_1(x) = e^x - 4(1 - x^2)$. Тогда $f'_1(x) = e^x + 8x$ и для $x \in [0.70, 0.72]$ верны оценки $0 < m = f'_1(0.70) \leq f'_1(x) \leq f'_1(0.72) = M$. Выберем в уравнении (4.26) $\alpha = 2/(m + M)$, возьмем $x^{(0)} = 0.7$ и будем вести итерации по формуле $x^{(n+1)} = x^{(n)} - \alpha[e^{x^{(n)}} - 4(1 - (x^{(n)})^2)]$. Выбранному α соответствует $q = (M - m)/(M + m) \approx 0.013$ и поэтому сходимость должна быть более быстрой, чем в примере 4.6. Действительно, уже первая итерация дает $x^{(1)} = 0.7034025118$, и так как $\frac{q}{1-q} |x^{(1)} - x^{(0)}| \approx 5 \cdot 10^{-5}$, то итерации следует прекратить и считать $\bar{x} = 0.7034 \pm 0.0001$. ▲

§ 4.5. Обусловленность метода простой итерации

В § 4.4 метод простой итерации был рассмотрен при идеальном предположении о возможности точного вычисления значений функции $\varphi(x)$. В действительности же вычисления на компьютере дают приближенные значения $\varphi^*(x)$. Поэтому вместо последовательности $x^{(n)}$, удовлетворяющей равенству $x^{(n+1)} = \varphi(x^{(n)})$, получается последовательность $\tilde{x}^{(n)}$, для которой

$$\tilde{x}^{(n+1)} = \varphi^*(\tilde{x}^{(n)}). \quad (4.27)$$

Известно, что метод простой итерации и многие другие итерационные методы устойчивы к ошибке, допущенной на одной из итераций. Такая ошибка эквивалентна некоторому ухудшению очередного приближения; если она не вывела приближение за пределы области сходимости, то итерационная последовательность по-прежнему будет сходиться к решению \bar{x} , а внесенная погрешность — затухать. Поэтому о таких итерационных методах говорят, что они обладают *свойством самоисправляемости*.

Однако погрешности допускаются не на одной, а на всех итерациях и совокупное их влияние несколько иное.

1. Обусловленность задачи. Прежде чем сформулировать результат о поведении метода простой итерации при наличии погрешности в вычислении функции φ , отметим, что преобразование уравнения $f(x) = 0$ к виду $x = \varphi(x)$ изменяет обусловленность задачи. Запишем это уравнение в виде $\tilde{f}(x) = 0$, где $\tilde{f}(x) = x - \varphi(x)$, и воспользуемся результатами § 4.2. Заметим, что $\bar{\Delta}(\tilde{f}^*) = \bar{\Delta}(\varphi^*)$, поскольку в действительности при-

ближенно вычисляется только функция ϕ . Поэтому оценка (4.11) в данном случае выглядит так:

$$\bar{\Delta}(\bar{x}^*) \leq \bar{\varepsilon} \approx v\bar{\Delta}(\phi^*). \quad (4.28)$$

Здесь $v = 1/|1 - \phi'(\bar{x})|$ — абсолютное число обусловленности корня \bar{x} . Грубо оценивая при выполнении условия $|\phi'| \leq q < 1$ величину v числом $\bar{v} = 1/(1 - q)$, приходим к оценке

$$\bar{\Delta}(\bar{x}^*) \lesssim \frac{\bar{\Delta}(\phi^*)}{1 - q}. \quad (4.29)$$

Когда $\phi' \leq 0$ или $\phi' \approx 0$, ее можно уточнить. Действительно, как нетрудно установить, $v \lesssim 1$, если $\phi' \lesssim 0$; следовательно,

$$\bar{\Delta}(\bar{x}^*) \lesssim \bar{\Delta}(\phi^*). \quad (4.30)$$

Заметим, что в оценках (4.28)–(4.30) величины $\bar{\Delta}(\bar{x}^*)$ и $\bar{\Delta}(\phi^*)$ можно заменить на $\bar{\delta}(\bar{x}^*)$ и $\bar{\delta}(\phi^*)$. Чтобы убедиться в этом, достаточно разделить левую и правую части оценок на величины $|\bar{x}|$ и $|\phi(\bar{x})|$, которые равны между собой. Например, оценка (4.28) преобразуется к виду

$$\bar{\delta}(\bar{x}^*) \lesssim v\bar{\delta}(\phi^*) \quad (4.31)$$

и, следовательно, абсолютное и относительное числа обусловленности здесь совпадают.

Сделаем некоторые выводы. Задача вычисления корня \bar{x} уравнения $x = \phi(x)$ плохо обусловлена, если $\phi'(\bar{x}) \approx 1$. В этом случае следует ожидать, что количество верных цифр корня \bar{x} по сравнению с количеством верных цифр в вычисляемых значениях $\phi^*(x)$ должно быть меньше примерно на $N = \lg v$ цифр. Для радиуса интервала неопределенности $\bar{\varepsilon} \approx v\bar{\Delta}(\phi^*)$ корня \bar{x} при $|\phi'| \leq q < 1$ справедлива оценка $\bar{\varepsilon} \leq \bar{\varepsilon}^* = \frac{\bar{\Delta}(\phi^*)}{1 - q}$. Когда $-1 < \phi' \leq 0$, уточненная оценка такова. $\bar{\varepsilon} \lesssim \bar{\varepsilon}^* = \bar{\Delta}(\phi^*)$; здесь потери верных цифр быть не должно.

Пример 4.9. Для уравнения $x = \phi(x)$ при $\phi(x) = 0.9999x + 10^{-4}/\sqrt{2}$ имеем $\phi'(x) = 0.9999$ и, следовательно, $v = 10^4$. Поэтому при решении этого уравнения методом простой итерации на компьютере будет потеряно примерно четыре значащих цифры. Вычисления на 6-разрядном десятичном компьютере (или на близком ему по точности компьютере типа IBM PC) могут дать всего лишь две верные значащие цифры. Это вполне согласуется с результатом, полученным в примере 4.6. ▲

2. Чувствительность метода простых итераций к погрешности вычислений. Сформулируем основной результат данного параграфа.

Теорема 4.5. Пусть выполнены условия теоремы 4.2 и для всех $x \in (\bar{x} - \sigma, \bar{x} + \sigma)$ имеет место неравенство $|\varphi(x) - \varphi^*(x)| \leq \bar{\Delta}(\varphi^*)$. Предположим также, что $\bar{\varepsilon}^* = \frac{\bar{\Delta}(\varphi^*)}{1-q} < \sigma$ (т.е. величина $\bar{\Delta}(\varphi^*)$ достаточно мала).

Если вычисления по формулам (4.16) и (4.27) начинаются с одного начального приближения $x^{(0)} = \tilde{x}^{(0)} \in (\bar{x} - \sigma_1, \bar{x} + \sigma_1)$, где $\sigma_1 =$

$= \min \left\{ \sigma, \frac{\sigma - \bar{\varepsilon}^*}{q} \right\}$, то последовательность $\tilde{x}^{(n)}$ не выходит за пределы

σ -окрестности корня \bar{x} и для всех $n \geq 1$ справедливы следующие оценки погрешности:

$$|\tilde{x}^{(n)} - x^{(n)}| \leq \bar{\varepsilon}^*, \quad (4.32)$$

$$|\tilde{x}^{(n)} - \bar{x}| \leq q^n |\tilde{x}^{(0)} - \bar{x}| + \bar{\varepsilon}^*, \quad (4.33)$$

$$|\tilde{x}^{(n)} - \bar{x}| \leq C |\tilde{x}^{(n)} - \tilde{x}^{(n-1)}| + \bar{\varepsilon}^*; \quad (4.34)$$

здесь $C = \frac{q}{1-q}$.

Замечание 1. При $\varphi' \leq 0$ в неравенствах (4.32)–(4.34) можно положить $\bar{\varepsilon}^* \approx \bar{\Delta}(\varphi^*)$, $C \approx 1$.

Замечание 2. При достаточно больших n в оценках (4.32)–(4.34) величину $\bar{\varepsilon}^*$ можно считать приближенно равной радиусу интервала неопределенности $\bar{\varepsilon} = \frac{\bar{\Delta}(\varphi^*)}{1 - \varphi'(\bar{x})}$.

Доказательство. Докажем по индукции, что для всех $n \geq 0$ справедливы неравенства (4.32), (4.33) и $\tilde{x}^{(n)} \in (\bar{x} - \sigma, \bar{x} + \sigma)$. Очевидно, что при $n = 0$ это верно.

Предположим, что доказываемое утверждение справедливо при некотором $n \geq 0$. Вычитая из равенства (4.27) равенство (4.16), получаем

$$\tilde{x}^{(n+1)} - x^{(n+1)} = \varphi'(\tilde{\xi}^{(n)}) (\tilde{x}^{(n)} - x^{(n)}) + \tilde{\Delta}^{(n)}, \quad \tilde{\xi}^{(n)} \in (\bar{x} - \sigma, \bar{x} + \sigma),$$

где $\tilde{\Delta}^{(n)} = \varphi^*(\tilde{x}^{(n)}) - \varphi(\tilde{x}^{(n)})$, $|\tilde{\Delta}^{(n)}| \leq \bar{\Delta}(\varphi^*)$.

Как следствие полученного равенства и сделанных предположений имеем

$$|\tilde{x}^{(n+1)} - x^{(n+1)}| \leq q |\tilde{x}^{(n)} - x^{(n)}| + |\tilde{\Delta}^{(n)}| \leq q \bar{\varepsilon}^* + \bar{\Delta}(\varphi^*) = \bar{\varepsilon}^*.$$

Объединяя эту оценку с оценкой (4.19), получаем

$$|\tilde{x}^{(n+1)} - \bar{x}| \leq q^{n+1} |\tilde{x}^{(0)} - \bar{x}| + \bar{\varepsilon}^*.$$

Учитывая, что $|\tilde{x}^{(0)} - \bar{x}| < \sigma_1$, имеем

$$|\tilde{x}^{(n+1)} - \bar{x}| < q\sigma_1 + \bar{\varepsilon}^* \leq \frac{q(\sigma - \bar{\varepsilon}^*)}{q} + \bar{\varepsilon}^* = \sigma.$$

Нужное утверждение доказано для номера, равного $n + 1$, а следовательно, и для всех $n \geq 0$.

Вычитая из равенства (4.27) равенство (4.17), получаем

$$\tilde{x}^{(n+1)} - \bar{x} = \alpha^{(n)}(\tilde{x}^{(n)} - \bar{x}) + \tilde{\Delta}^{(n)}, \quad \alpha^{(n)} = \varphi'(\xi^{(n)}), \quad \xi^{(n)} \in (\bar{x} - \sigma, \bar{x} + \sigma).$$

Таким образом,

$$\tilde{x}^{(n+1)} - \bar{x} = \frac{\alpha^{(n)}}{1 - \alpha^{(n)}} (\tilde{x}^{(n)} - \tilde{x}^{(n+1)}) + \frac{\tilde{\Delta}^{(n)}}{1 - \alpha^{(n)}}.$$

Из полученного равенства вытекает оценка (4.34). ■

Итак, метод простой итерации обладает такой же чувствительностью к погрешностям в вычислении функции φ , что и задача вычисления корня уравнения $x = \varphi(x)$. Метод позволяет получить приближенное значение корня \bar{x} с точностью, примерно совпадающей с радиусом $\bar{\varepsilon}$ интервала неопределенности корня. Критерий (4.23) окончания итераций применим, если $\bar{\varepsilon} \ll \varepsilon$ (радиус интервала неопределенности много меньше заданной точности); это следует из неравенства (4.34) и замечания 2.

Входящую в соотношения (4.29), (4.30) величину $\bar{\Delta}(\varphi^*)$ в общем случае оценить сверху достаточно сложно. Оценка снизу очевидна: $\bar{\Delta}(\varphi^*) \gtrsim |\bar{x}| \varepsilon_m$. В благоприятной ситуации, когда вычисления ведутся по простым формулам, можно надеяться на то, что $\bar{\Delta}(\varphi^*)$ окажется величиной порядка $|\bar{x}| \varepsilon_m$.

§ 4.6. Метод Ньютона

Знаменитый метод Ньютона является одним из наиболее эффективных методов решения самых разных нелинейных задач. Расчетную формулу метода можно получить, использовав различные подходы. Рассмотрим два из них.

1. Метод касательных. Выведем расчетную формулу метода для решения нелинейного уравнения (4.1) из простых геометрических соображений. Соответствующая иллюстрация приведена на рис. 4.9.

Пусть $x^{(0)}$ — заданное начальное приближение к корню \bar{x} . В точке $M^{(0)}$ с координатами $(x^{(0)}, f(x^{(0)}))$ проведем касательную к графику функции $y = f(x)$ и за новое приближение $x^{(1)}$ примем абсциссу точки пересечения этой касательной с осью Ox . Аналогично, за приближение $x^{(2)}$ примем абсциссу точки пересечения с осью Ox касательной, проведенной к графику в точке $M^{(1)}$ с координатами $(x^{(1)}, f(x^{(1)}))$. Продолжая этот процесс далее, получим последовательность $x^{(0)}, x^{(1)}, x^{(2)}, \dots, x^{(n)}, \dots$ приближений к корню \bar{x} .

Напомним, что уравнение касательной, проведенной к графику функции $y = f(x)$ в точке $(x^{(n)}, f(x^{(n)}))$, имеет вид

$$y = f(x^{(n)}) + f'(x^{(n)})(x - x^{(n)}). \quad (4.35)$$

Пусть $f'(x^{(n)}) \neq 0$. Полагая в равенстве (4.35) $y = 0$, замечаем, что абсцисса $x^{(n+1)}$ точки пересечения касательной с осью Ox удовлетворяет равенству

$$0 = f(x^{(n)}) + f'(x^{(n)})(x^{(n+1)} - x^{(n)}). \quad (4.36)$$

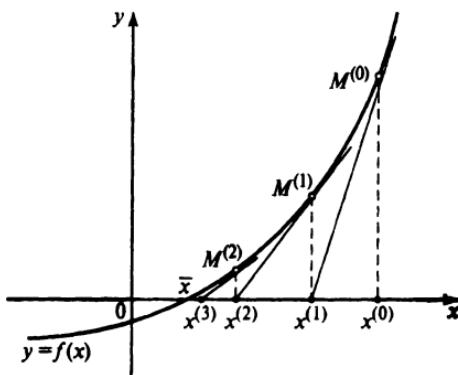


Рис. 4.9

Выражая из него $x^{(n+1)}$, получаем расчетную формулу *метода Ньютона*:

$$x^{(n+1)} = x^{(n)} - \frac{f(x^{(n)})}{f'(x^{(n)})}, \quad n \geq 0. \quad (4.37)$$

Благодаря такой геометрической интерпретации этот метод часто называют *методом касательных*.

2. Метод линеаризации. С более общих позиций метод Ньютона можно рассматривать как итерационный метод, использующий специальную линеаризацию задачи и позволяющий свести решение исходного нелинейного уравнения к решению последовательности линейных уравнений.

Пусть приближение $x^{(n)}$ уже получено. Представим функцию в окрестности точки $x^{(n)}$ по формуле Тейлора:

$$f(x) = f(x^{(n)}) + f'(x^{(n)})(x - x^{(n)}) + \frac{f''(\xi)}{2} (x - x^{(n)})^2. \quad (4.38)$$

Здесь ξ — некоторая точка, расположенная между x и $x^{(n)}$. Заменяя в уравнении $f(x) = 0$ функцию $f(x)$ главной линейной частью разложений (4.38), получаем линейное уравнение

$$f(x^{(n)}) + f'(x^{(n)})(x - x^{(n)}) = 0. \quad (4.39)$$

Принимая решение уравнения (4.39) за новое приближение $x^{(n+1)}$, приходим к формуле (4.37).

3. Основная теорема о сходимости метода Ньютона.

Теорема 4.6. Пусть \bar{x} — простой корень уравнения $f(x) = 0$, в некоторой окрестности которого функция f дважды непрерывно дифференцируема. Тогда найдется такая малая σ -окрестность корня \bar{x} , что при произвольном выборе начального приближения $x^{(0)}$ из этой окрестности итерационная последовательность метода Ньютона не выходит за пределы окрестности и справедлива оценка

$$|x^{(n+1)} - \bar{x}| \leq \sigma^{-1} |x^{(n)} - \bar{x}|^2, \quad n \geq 0, \quad (4.40)$$

означающая, что метод сходится с квадратичной скоростью.

Следствием оценки (4.40) является априорная оценка

$$|x^{(n)} - \bar{x}| \leq \sigma q^{2^n}, \quad n \geq 0, \quad (4.41)$$

в которой $q = \sigma^{-1} |x^{(0)} - \bar{x}|$.

Доказательство. Так как $f'(\bar{x}) \neq 0$ (по определению простого корня), то в силу непрерывности функций f' и f'' найдется δ_0 -окрестность корня, в которой при некоторых постоянных α и β выполнены неравенства $0 < \alpha \leq |f'(x)|$, $|f''(x)| \leq \beta$.

Пусть $x^{(n)} \in (\bar{x} - \sigma, \bar{x} + \sigma)$, где $\sigma = \min\{\delta_0, 2\alpha/\beta\}$.

Подставляя $x = \bar{x}$ в (4.38), получаем равенство

$$0 = f(x^{(n)}) + f'(x^{(n)})(\bar{x} - x^{(n)}) + \frac{f'(\xi)}{2} (\bar{x} - x^{(n)})^2,$$

в котором $\xi \in (\bar{x} - \sigma, \bar{x} + \sigma)$. Вычитая из него равенство (4.36), имеем

$$f'(x^{(n)})(x^{(n+1)} - \bar{x}) = \frac{f''(\xi)}{2} (\bar{x} - x^{(n)})^2.$$

Тогда, приравнивая модули обеих частей этого равенства и используя условия ограниченности $|f'(x)|$ и $|f''(x)|$, приходим к неравенству

$$\alpha |x^{(n+1)} - \bar{x}| \leq \frac{\beta}{2} |x^{(n)} - \bar{x}|^2,$$

откуда следует справедливость оценки (4.40). Доказательство теоремы завершается применением леммы 4.2. ■

Таким образом, при выборе начального приближения из достаточно малой окрестности корня метод Ньютона сходится квадратично. Это означает, грубо говоря, что на каждой итерации число верных цифр приближения примерно удваивается.

Приведенные в теореме 4.6 оценки погрешности являются априорными (см. § 3.3) и их использование в практике вычислений для количественной оценки погрешности неэффективно или чаще всего невозможно.

4. Критерий окончания. На практике предпочтительнее использование простой апостериорной оценки

$$|x^{(n)} - \bar{x}| \leq |x^{(n)} - x^{(n-1)}|, \quad (4.42)$$

справедливость которой обосновывается следующим утверждением.

Теорема 4.7. Пусть выполнены условия теоремы 4.6 и $|x^{(0)} - \bar{x}| < \sigma/2$. Тогда для всех $n \geq 1$ верна оценка (4.42).

Доказательство. Из оценки (4.41) следует, что

$$|x^{(n-1)} - \bar{x}| \leq \sigma q^{2^{n-1}} \leq \sigma q = |x^{(0)} - \bar{x}| < \frac{\sigma}{2}.$$

Поэтому, применив неравенство (4.40), получаем цепочку неравенств

$$\begin{aligned} 2|x^{(n)} - \bar{x}| &\leq 2\sigma^{-1}|x^{(n-1)} - \bar{x}|^2 \leq |x^{(n-1)} - \bar{x}| \leq \\ &\leq |x^{(n-1)} - x^{(n)}| + |x^{(n)} - \bar{x}|, \end{aligned}$$

из которой вытекает оценка (4.42). ■

Наличие оценки (4.42) позволяет сформулировать следующий практический критерий окончания итераций метода Ньютона. При заданной точности $\varepsilon > 0$ вычисления нужно вести до тех пор, пока не окажется выполненным неравенство

$$|x^{(n)} - x^{(n-1)}| < \varepsilon. \quad (4.43)$$

Пример 4.8. Используя метод Ньютона, найдем с точностью $\varepsilon = 10^{-6}$ положительный корень уравнения $4(1 - x^2) - e^x = 0$.

В примере 4.2 корень был локализован на отрезке $[0, 1]$. Для $f(x) = 4(1 - x^2) - e^x$ имеем $f'(x) = -8x - e^x$. Очевидно, что $f'(\bar{x}) \neq 0$, т.е. \bar{x} — простой корень. Возьмем начальное приближение $x^{(0)} = 0.5$ и будем выполнять итерации метода Ньютона по формуле

$$x^{(n+1)} = x^{(n)} + \frac{4(1 - (x^{(n)})^2) - e^{x^{(n)}}}{8x^{(n)} + e^{x^{(n)}}}$$

Результаты первых итераций с десятью знаками мантиссы приведены в табл. 4.4.

Таблица 4.4

n	$x^{(n)}$	$ x^{(n)} - x^{(n-1)} $
0	0.5000000000	—
1	0.7392185177	$2.4 \cdot 10^{-1}$
2	0.7042444088	$3.5 \cdot 10^{-2}$
3	0.7034399951	$8.0 \cdot 10^{-4}$
4	0.7034395712	$4.3 \cdot 10^{-7}$

При $n = 4$ вычисления следует прекратить и после округления получим $\bar{x} = 0.703440 \pm 0.000001$.

Сравнение результатов итераций со значением \bar{x} показывает, что приближения $x^{(1)}, x^{(2)}, x^{(3)}$ содержат одну, три, шесть верных значащих цифр соответственно. Можно показать, что приближение $x^{(4)}$ содержит десять верных цифр. Это подтверждает отмеченный ранее общий факт: при каждой итерации метода Ньютона число верных значащих цифр примерно удваивается. ▲

Пример 4.9. Используя метод Ньютона, укажем итерационный процесс вычисления $\sqrt[p]{a}$, где $a > 0$, p — натуральное число.

По определению, $x = \sqrt[p]{a}$ — это неотрицательная величина, удовлетворяющая равенству $x^p = a$. Таким образом, задача сводится к вычислению положительного корня уравнения $f(x) = 0$, где $f(x) = x^p - a$. Итерационная формула метода Ньютона примет вид

$$x^{(n+1)} = x^{(n)} - \frac{(x^{(n)})^p - a}{p(x^{(n)})^{p-1}} = \frac{p-1}{p} x^{(n)} + \frac{a}{p(x^{(n)})^{p-1}}. \quad (4.44)$$

При $p = 2$ эта формула уже была получена в примере 3.17. ▲

5. Связь с методом простой итерации. Метод Ньютона можно рассматривать как один из вариантов метода простой итерации, связанный со специальным преобразованием уравнения $f(x) = 0$ к виду

$$x = \varphi_N(x), \quad (4.45)$$

где $\varphi_N(x) = x - f(x)/f'(x)$.

Действительно, итерационная формула метода простой итерации $x^{(n+1)} = \varphi_N(x^{(n)})$ совпадает с (4.37).

Если исходить из этого с учетом оценки (4.19), можно сделать вывод о том, что метод Ньютона сходится только линейно. Однако заметим, что $\varphi'_N(x) = f(x) \frac{f''(x)}{(f'(x))^2}$. Так как $f(\bar{x}) = 0$, то $\varphi'_N(\bar{x}) = 0$ и величина $\alpha^{(n+1)} = |\varphi'(\xi^{(n)})|$, определяющая в силу равенства (4.20) коэффициент сжатия погрешности, стремится к нулю при $n \rightarrow \infty$. Скорость сходимости возрастает по мере приближения к корню, отсюда и ее сверхлинейный характер.

В качестве аналога теоремы 4.4 для метода Ньютона приведем следующий результат.

Теорема 4.8. Пусть $[a, b]$ — отрезок локализации простого корня \bar{x} уравнения (4.1). Предположим, что на этом отрезке функция f дважды непрерывно дифференцируема, а ее производные $f'(x)$ и $f''(x)$ знакопостоянны. Пусть $x^{(0)} \in [a, b]$ — произвольное начальное приближение, и в случае $f(x^{(0)})f''(x^{(0)}) < 0$ выполнено дополнительное условие $x^{(1)} \in [a, b]$ (первое приближение не выходит за пределы отрезка локализации).

Тогда начиная с $n = 1$ итерационная последовательность метода Ньютона $x^{(n)}$ сходится к \bar{x} монотонно с той стороны отрезка $[a, b]$, где $f(x)f''(x) > 0$. ■

Иллюстрацией монотонного характера сходимости может служить рис. 4.9.

Следствие. Пусть уравнение $f(x)$ имеет корень \bar{x} , функция $f(x)$ дважды непрерывно дифференцируема на всей числовой оси, а ее производные f' и f'' знакопостоянны. Тогда метод Ньютона сходится при любом начальном приближении $x^{(0)}$ (т.е. является глобально сходящимся), причем начиная с $n = 1$ последовательность сходится монотонно с той стороны от корня, где $f(x)f''(x) > 0$.

6. Трудности использования. Простота, логическая стройность и высокая скорость сходимости делают метод Ньютона чрезвычайно привлекательным. Однако для его практического применения нужно преодолеть две существенные трудности. Одна из них состоит в необходимости вычисления производной $f'(x)$. Часто бывает невозможно найти аналитическое выражение для $f'(x)$, а определить приближенное значение с высокой точностью очень трудно. Иногда вычисление $f'(x)$ — вполне реальная, но весьма дорогостоящая операция. В этих случаях приходится модифицировать метод, избегая непосредственного вычисления производной. Некоторые из таких модификаций приведены в § 4.7.

Более существенно то, что метод Ньютона обладает, вообще говоря, только локальной сходимостью. Это означает, что областью его сходимости является некоторая малая σ -окрестность корня \bar{x} и для гарантии сходимости необходимо выбирать хорошее начальное приближение, попадающее в эту σ -окрестность. Неудачный выбор начального приближения может дать расходящуюся последовательность (рис. 4.10) и даже привести к преждевременному прекращению вычислений (если на очередной итерации $f'(x^{(n)}) \approx 0$). Для преодоления этой трудности часто используют метод Ньютона в сочетании с каким-либо медленно, но гарантированно сходящимся методом типа бисекции. Такие гибридные алгоритмы находят в последнее время широкое практическое применение.

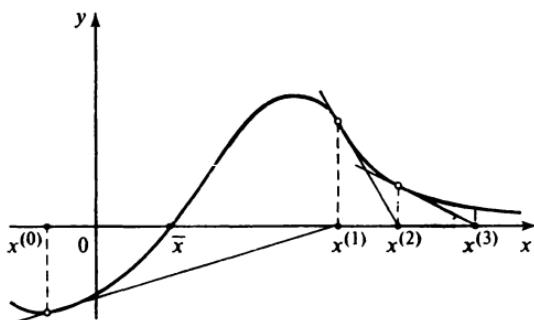


Рис. 4.10

7. Влияние погрешности вычислений. Пусть \bar{x} — простой корень. Если метод Ньютона рассматривать как вариант метода простой итерации, связанный с преобразованием уравнения $f(x) = 0$ к виду (4.45), то можно воспользоваться результатами § 4.4 и 4.5.

Поскольку $\phi'_N(\bar{x}) = 0$, то справедливо неравенство (4.30) и радиус интервала неопределенности $\bar{\varepsilon}_N$ корня \bar{x} преобразованного уравнения равен примерно $\bar{\Delta}(\phi_N^*)$. В оценках (4.32)–(4.34) можно считать $\bar{\varepsilon}^* \approx \bar{\Delta}(\phi_N^*)$ и $C \approx 1$. Для того чтобы сделать окончательные выводы, остается оценить $\bar{\Delta}(\phi_N^*)$.

Пусть f^* и $(f')^*$ — приближенные значения функций f и f' , вычисляемые в малой окрестности корня \bar{x} . Будем считать, что производная f' вычисляется хотя бы с точностью до одной-двух верных значащих цифр. В противном случае (особенно, если неверно определяется знак f') из геометрического смысла метода легко понять, что его применять не следует. В силу предложения 2.5 в малой окрестности корня имеем

$$\bar{\Delta}\left(\frac{f^*}{(f')^*}\right) \approx \frac{\bar{\Delta}(f^*)}{|f'|} + \frac{\bar{\Delta}((f')^*)|f|}{|f'|^2} \approx \frac{\bar{\Delta}(f^*)}{|f'(\bar{x})|} = \bar{\varepsilon}.$$

Так как $\phi_N^*(x) = x \ominus \frac{f^*(x)}{(f')^*(x)} \approx \bar{x}$ то, учитывая погрешность, приближенно равную $|\bar{x}| \varepsilon_m$, вносимую в ϕ_N^* вследствие выполнения операции вычитания, получаем $\bar{\varepsilon}_N = \bar{\Delta}(\phi_N^*) \approx \bar{\varepsilon} + |\bar{x}| \varepsilon_m$.

Таким образом, преобразование уравнения $f(x) = 0$ к виду (4.45) практически не меняет обусловленность и радиус $\bar{\varepsilon}$ интервала неопределенности корня \bar{x} . Итерационный процесс Ньютона дает возможность вычислить решение с точностью $\varepsilon \geq \bar{\varepsilon}$. Отметим, тем не менее, что эти достоинства метода Ньютона реализуются, вообще говоря, только в малой окрестности корня.

§ 4.7. Модификации метода Ньютона

В § 4.6 в качестве недостатка метода Ньютона была отмечена необходимость вычисления значения производной $f'(x)$ на каждой итерации. Рассмотрим некоторые модификации метода Ньютона, свободные от этого недостатка. Заметим, что, по существу, излагаемые в этом параграфе итерационные методы решения нелинейного уравнения на каждой итерации используют некоторую процедуру его линеаризации,

т.е. исходное нелинейное уравнение заменяется приближенно более простым линейным уравнением.

1. Упрощенный метод Ньютона. Если производная $f'(x)$ непрерывна, то ее значение вблизи простого корня \bar{x} почти постоянно. Поэтому можно попытаться вычислить f' лишь однажды в точке $x^{(0)}$, а затем заменить в формуле (4.37) значение $f'(x^{(n)})$ постоянной $f'(x^{(0)})$. В результате получим расчетную формулу *упрощенного метода Ньютона*:

$$x^{(n+1)} = x^{(n)} - \frac{f(x^{(n)})}{f'(x^{(0)})}, \quad n \geq 0. \quad (4.46)$$

Геометрическая иллюстрация метода приведена на рис. 4.11. В точке $(x^{(0)}, f(x^{(0)}))$ к графику функции $y = f(x)$ проводится касательная l_0 и за приближение $x^{(1)}$ принимается абсцисса точки пересечения этой касательной с осью Ox (как в методе Ньютона). Каждое следующее приближение $x^{(n+1)}$ получается здесь как абсцисса точки пересечения с осью Ox прямой, проходящей через точку $M^{(n)}$ с координатами $(x^{(n)}, f(x^{(n)}))$ и параллельной касательной l_0 .

Упрощение вычислений по сравнению с методом Ньютона достигается здесь ценой резкого падения скорости сходимости. Сходимость этого метода является уже не квадратичной, а линейной.

Метод (4.46) можно рассматривать как метод простой итерации с итерационной функцией $\varphi(x) = x - \frac{f(x)}{f'(x^{(0)})}$. Так как $\varphi'(x) = 1 - \frac{f'(x)}{f'(x^{(0)})}$, то для знаменателя q соответствующей геометрической прогрессии имеем $q \approx \left| 1 - \frac{f'(\bar{x})}{f'(x^{(0)})} \right|$. Следовательно, скорость сходимости тем выше, чем ближе начальное приближение $x^{(0)}$ к решению \bar{x} .

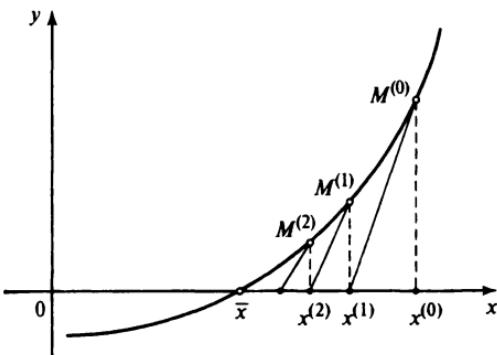


Рис. 4.11

2. Метод ложного положения. В основе этой и следующих двух модификаций метода Ньютона лежит приближенное равенство

$$f'(x^{(n)}) \approx \frac{f(z^{(n)}) - f(x^{(n)})}{z^{(n)} - x^{(n)}}. \quad (4.47)$$

Оно верно при условии $z^{(n)} \approx x^{(n)}$ и следует из определения производной:

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

Пусть c — фиксированная точка, расположенная в окрестности простого корня \bar{x} . Заменим в расчетной формуле метода Ньютона (4.37) производную $f'(x^{(n)})$ правой частью приближенного равенства (4.47), полагая $z^{(n)} = c$. В результате придем к расчетной формуле *метода ложного положения*:

$$x^{(n+1)} = x^{(n)} - \frac{c - x^{(n)}}{f(c) - f(x^{(n)})} f(x^{(n)}), \quad n \geq 0. \quad (4.48)$$

Геометрическая иллюстрация метода приведена на рис. 4.12. Очередное приближение $x^{(n+1)}$ получается здесь как абсцисса точки пересечения с осью Ox прямой, проведенной через расположенные на графике функции $y = f(x)$ точки M и $M^{(n)}$ с координатами $(c, f(c))$ и $(x^{(n)}, f(x^{(n)}))$.

Метод (4.48) обладает только линейной сходимостью. Его можно рассматривать как метод простой итерации с итерационной функцией $\varphi(x) = x - \frac{c - x}{f(c) - f(x)} f(x)$. Так как скорость сходимости определяется

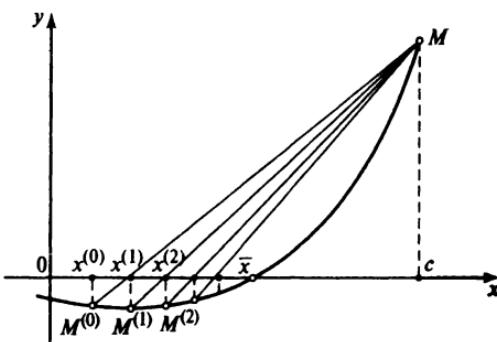


Рис. 4.12

вблизи корня величиной $q \approx |\varphi'(\bar{x})| = \left| 1 - \frac{(c-x)f'(\bar{x})}{f(c)-f(\bar{x})} \right|$, то она тем выше, чем ближе окажется выбранная точка c к \bar{x} .

3. Метод секущих. Замена в формуле метода Ньютона производной $f'(x^{(n)})$ приближением $\frac{f(x^{(n-1)}) - f(x^{(n)})}{x^{(n-1)} - x^{(n)}}$ приводит к расчетной формуле *метода секущих*:

$$x^{(n+1)} = x^{(n)} - \frac{x^{(n-1)} - x^{(n)}}{f(x^{(n-1)}) - f(x^{(n)})} f(x^{(n)}), \quad n \geq 0. \quad (4.49)$$

Заметим, что этот метод двухшаговый, так как для нахождения очередного приближения $x^{(n+1)}$ требуется знание двух предыдущих приближений $x^{(n)}$ и $x^{(n-1)}$. В частности, для того чтобы начать вычисления, необходимо задать два начальных приближения $x^{(0)}$ и $x^{(1)}$. Все рассмотренные ранее методы требовали для вычисления $x^{(n+1)}$ только знание $x^{(n)}$, т.е. были одношаговыми.

На рис. 4.13 приведена геометрическая иллюстрация метода. Очередное приближение $x^{(n+1)}$ получается здесь как абсцисса точки пересечения с осью Ox секущей, соединяющей точки $M^{(n-1)}$ и $M^{(n)}$ графика функции $f(x)$ с координатами $(x^{(n-1)}, f(x^{(n-1)}))$ и $(x^{(n)}, f(x^{(n)}))$.

Примечательно то, что эта модификация метода Ньютона сохраняет свойство сверхлинейной сходимости, если вычисляется простой корень \bar{x} . Точнее, верно следующее утверждение.

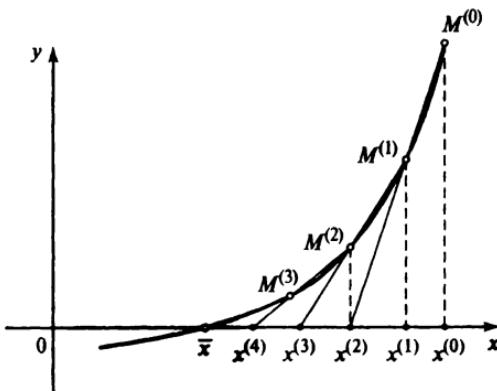


Рис. 4.13

Теорема 4.9. Пусть \bar{x} — простой корень уравнения $f(x) = 0$, в некоторой окрестности которого функция f дважды непрерывно дифференцируема, причем $f''(\bar{x}) \neq 0$. Тогда существует σ -окрестность корня \bar{x} такая, что при произвольном выборе приближений $x^{(0)}$ и $x^{(1)}$ из этой σ -окрестности метод секущих сходится с порядком $p = \frac{\sqrt{5} + 1}{2} \approx 1.618$, т.е. для $n \geq 1$ справедлива оценка

$$|x^{(n+1)} - \bar{x}| \leq c |x^{(n)} - \bar{x}|^p, \quad p = \frac{\sqrt{5} + 1}{2}. \blacksquare$$

Так как одна итерация метода секущих требует только одного нового вычисления значения функции f , а метод Ньютона — двух вычислений значений функций (f и f'), то трудоемкость двух итераций метода секущих приблизительно эквивалентна трудоемкости одной итерации по Ньютону. Две итерации метода секущих дают порядок $p^2 \approx 2.618 > 2$, поэтому его можно расценивать как более быстрый по сравнению с методом Ньютона.

К сожалению, метод обладает, вообще говоря, только локальной сходимостью. Он требует выбора двух близких к \bar{x} (в общем случае — очень близких) начальных приближений $x^{(0)}$ и $x^{(1)}$. Если эти приближения выбраны неудачно, то метод расходится (рис. 4.14).

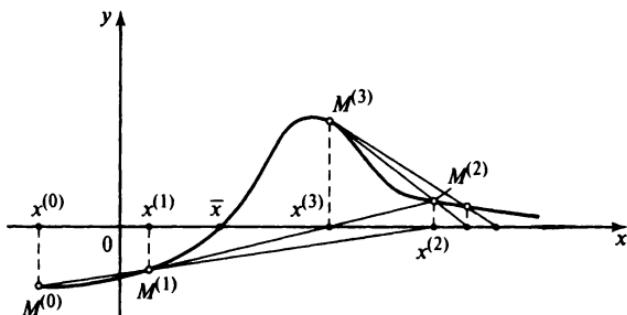


Рис. 4.14

4. Метод Стеффенсена. Итерационная формула метода Стеффенсена имеет вид

$$x^{(n+1)} = x^{(n)} - \frac{f(x^{(n)})}{f(x^{(n)} + f(x^{(n)})) - f(x^{(n)})} f(x^{(n)}), \quad n \geq 0. \quad (4.50)$$

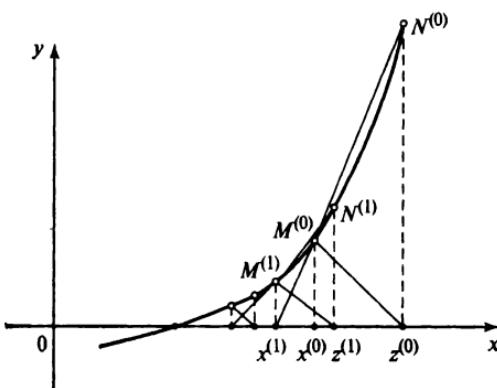


Рис. 4.15

Можно считать, что она получена в результате замены производной $f'(x^{(n)})$, входящей в расчетную формулу метода Ньютона, приближением (4.47), где $z^{(n)} = x^{(n)} + f(x^{(n)})$.

Метод Стеффенсена интересен тем, что он является одношаговым, не требует вычисления производной f' и в то же время, как и метод Ньютона, сходится квадратично, если корень \bar{x} — простой, функция f дважды непрерывно дифференцируема в окрестности корня, а начальное приближение $x^{(0)}$ выбрано близко к \bar{x} .

Геометрическая иллюстрация метода Стеффенсена приведена на рис. 4.15. Приближение $x^{(n+1)}$ получается как абсцисса точки пересечения с осью Ox секущей, проходящей через точки $M^{(n)}$ и $N^{(n)}$ с координатами $(x^{(n)}, f(x^{(n)}))$ и $(z^{(n)}, f(z^{(n)}))$. Значение $z^{(n)}$ отвечает абсциссе точки пересечения с осью Ox прямой $y = f(x^{(n)}) - (x - x^{(n)})$, проходящей через точку $M^{(n)}$ и параллельной прямой $y = -x$.

Несмотря на свойство квадратичной сходимости, метод Стеффенсена уступает методу секущих, поскольку требует большей вычислительной работы для достижения той же точности ε . Это связано с тем, что на каждой итерации рассматриваемого метода вычисление функции производится дважды, а в методе секущих лишь один раз.

5. Уточнение метода Ньютона для случая кратного корня. В принципе для вычисления корня уравнения $f(x) = 0$ кратности $m > 1$ можно использовать и стандартный метод Ньютона. Однако в этом случае скорость его сходимости является только линейной. Можно показать, что знаменатель q соответствующей геометрической прогрессии приближенно равен $1 - 1/m$.

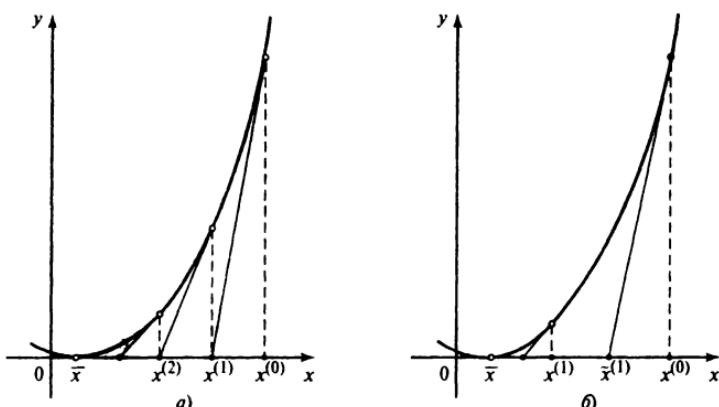


Рис. 4.16

Для того чтобы сохранить квадратичную скорость сходимости, метод Ньютона нужно модифицировать следующим образом:

$$x^{(n+1)} = x^{(n)} - m \frac{f(x^{(n)})}{f'(x^{(n)})}, \quad n \geq 0. \quad (4.51)$$

Можно показать (это достигается раскрытием неопределенностей с помощью формулы Тейлора), что при таком выборе итерационной функции $\phi(x) = x - m \frac{f(x)}{f'(x)}$ получим $\phi'(\bar{x}) = 0$, и сходимость снова окажется квадратичной.

На рис. 4.16, а, б проиллюстрировано поведение последовательных приближений стандартного метода Ньютона и его модификации (4.51) для случая отыскания корня кратности $m = 2$.

Для метода (4.51) значение $x^{(n+1)}$ получается следующим образом. В точке $M^{(n)}$ с координатами $(x^{(n)}, f(x^{(n)}))$ к графику функции проводится касательная. Пересечение ее с осью Ox дает вспомогательную точку $\tilde{x}^{(n+1)}$. Для получения точки $x^{(n+1)}$ нужно воспользоваться равенством $x^{(n+1)} - x^{(n)} = m(\tilde{x}^{(n+1)} - x^{(n)})$.

Пример 4.10. Применим методы, рассмотренные в данной главе, для вычисления положительного корня уравнения $4(1 - x^2) - e^x = 0$, считая известным, что $\bar{x} \in [0, 1]$ (см. пример 4.2).

Результаты вычислений приведены в табл. 4.5—4.8. В них для каждого приближения дается число верных знаков и требуемое число вычислений значений функции $f(x) = 4(1 - x^2) - e^x$ (для упрощенного метода

Таблица 4.5

Упрощенный метод Ньютона;
 $x^{(0)} = 0.5$

n	$x^{(n)}$	Число верных знаков	Число вычислений функций
0	0.5000000000	0	0
1	0.7392185177	1	2
2	0.6896366262	1	3
3	0.7081565866	2	4
4	0.7017501401	2	5
5	0.7040350503	3	6
6	0.7032284740	3	7
7	0.7035142540	4	8
8	0.7034131306	4	9
9	0.7034489297	4	10
10	0.7034362584	5	11

Таблица 4.7

Метод секущих; $x^{(0)} = 1, x^{(1)} = 0.5$

n	$x^{(n)}$	Число верных знаков	Число вычислений функций
0	1.0000000000	0	0
1	0.5000000000	0	0
2	0.6660226835	1	2
3	0.7092548653	2	3
4	0.7032943284	3	4
5	0.7034390197	6	5
6	0.7034395712	10	6

Таблица 4.6

Метод ложного положения;
 $c = 1, x^{(0)} = 0.5$

n	$x^{(n)}$	Число верных знаков	Число вычислений функций
0	0.5000000000	0	0
1	0.6660226835	1	2
2	0.6971284254	2	3
3	0.7023912700	2	4
4	0.7032658920	3	5
5	0.7034108088	4	6
6	0.7034348083	5	7
7	0.7034387825	6	8
8	0.7034394406	6	9
9	0.7034395495	7	10
10	0.7034395676	8	11

Таблица 4.8

Метод Стеффенсена; $x^{(0)} = 0.5$

n	$x^{(n)}$	Число верных знаков	Число вычислений функций
0	0.5000000000	0	0
1	0.6047702913	0	2
2	0.6731754719	1	4
3	0.6998880287	2	6
4	0.7033854199	3	8
5	0.7034395584	7	10
6	0.7034395712	10	12

Ньютона учтено вычисление $f'(x^{(0)})$). Вычисления выполнены с десятью знаками мантиссы.

Отметим, что выбор начальных приближений был довольно случайным (хотя и разумным). Тем не менее лучший результат показал метод секущих. Решение с десятью верными знаками мантиссы было получено после пяти итераций и для этого потребовалось лишь шесть вычислений функции. Хотя метод Ньютона (см. пример 4.8) при том же начальном приближении $x^{(0)}$ дает такое же значение \bar{x} всего после четырех итераций, для этого требуется восемь вычислений функции (четыре вычисления f и четыре вычисления f'). ▲

Таблица 4.9

Метод Ньютона	
n	$x^{(n)}$
0	1.0000000000
1	0.6666666667
2	0.4166666667
3	0.2442528736
4	0.1354180464
5	0.0720028071
6	0.0372524647
7	0.0189668249
8	0.0095725027

Таблица 4.10
Уточнение метода Ньютона
для случая $m = 2$

n	$x^{(n)}$
0	1.0000000000
1	0.3333333333
2	0.0476190476
3	0.0011074197
4	0.0000006128
5	0.0000000000

Пример 4.11. Применим метод Ньютона и его модификацию (4.51) для вычисления корня $\bar{x} = 0$ кратности $m = 2$ для уравнения $x^2 e^x = 0$. Возьмем $x^{(0)} = 1$. Результаты вычислений даны в табл. 4.9 и 4.10.

Как видно из табл. 4.9, погрешность метода Ньютона убывает довольно медленно и примерно соответствует геометрической прогрессии со знаменателем $q = 1/2$. В то же время пять итераций по формуле (4.51) при $m = 2$ дают значение решения с погрешностью, меньшей $\varepsilon = 10^{-10}$. ▲

6. Чувствительность к погрешностям. Рассмотренные в этом параграфе одношаговые методы можно интерпретировать как различные варианты метода простой итерации. Поэтому исследование их чувствительности к погрешностям сводится (аналогично тому, как это было сделано в § 4.6 для метода Ньютона) к использованию соответствующих результатов § 4.5. Например, можно убедиться в хорошей обусловленности модифицированного метода Ньютона и метода ложного положения.

Высокая скорость сходимости метода секущих делает его привлекательным для применения. Однако вызывает беспокойство тот факт, что в формулу (4.49) входит величина $\frac{f(x^{(n-1)}) - f(x^{(n)})}{x^{(n-1)} - x^{(n)}}$, аппроксимирующая производную. Вблизи корня, когда $x^{(n)} \approx x^{(n-1)}$ и $f(x^{(n)}) \approx f(x^{(n-1)}) \approx 0$, погрешность вычисления функции начинает существенно сказываться на точности этой величины, и метод секущих теряет устойчивость. Этим он существенно отличается от методов простой итерации и Ньютона, для которых приближения $x^{(n)}$ равномерно устойчивы к погрешности независимо от числа итераций n .

Тем не менее при грамотном использовании метод секущих дает возможность получить почти tanto же верных значащих цифр корня \bar{x} , сколько вообще позволяет обусловленность задачи (см. § 4.2). Воз-

можная (но не обязательная) потеря точности составляет одну-две верные цифры. Необходимо лишь прервать итерационный процесс в тот момент, когда приближения окажутся в опасной близости к интервалу неопределенности. Один из способов заметить этот момент состоит в использовании правила Гарвика (см. § 4.2).

Отметим, что при попадании очередного приближения в малую окрестность решения теряет устойчивость и метод Стеффенсена.

§ 4.8. Дополнительные замечания

1. Метод обратной квадратичной интерполяции. Пусть приближения $x^{(n-2)}$, $x^{(n-1)}$, $x^{(n)}$ уже найдены, а соответствующие значения функции f различны (последнее предположение заведомо будет выполнено, если функция f строго монотонна). Тогда строят такой квадратичный многочлен $P_2(y)$ от переменной y , что $x^{(i)} = P_2(f(x^{(i)}))$ для $i = n-2, n-1, n$. За очередное приближение к решению принимается $x^{(n+1)} = P_2(0)$.

Этот трехшаговый метод обладает локальной сходимостью с порядком $p \approx 1.839$; для начала его работы требуется задание трех хороших начальных приближений.

2. Методы с порядком сходимости выше второго. Приведем примеры одношаговых методов, обладающих кубической сходимостью. Таковым является итерационный метод¹

$$x^{(n+1)} = x^{(n)} - u(x^{(n)}) \left[1 + \frac{1}{2} v(x^{(n)}) \right], \quad n \geq 0.$$

Здесь

$$u(x) = \frac{f(x)}{f'(x)}, \quad v(x) = \frac{f(x)f''(x)}{[f'(x)]^2}.$$

Другим примером метода, обладающего кубической скоростью сходимости, является *метод Галлея*²

$$x^{(n+1)} = x^{(n)} - u(x^{(n)}) \frac{1}{1 - \frac{1}{2} v(x^{(n)})}, \quad n \geq 0.$$

¹ Этот метод часто называют итерационным методом Эйлера.

² Эдмунд Галлей (1656—1742) — английский астроном, предсказавший периодическое (с периодом 75 лет) появление кометы, впоследствии названной его именем.

Существуют и другие методы с порядком сходимости выше второго, теоретически очень быстро сходящиеся вблизи корня. Однако они редко используются на практике, так как их преимущество в скорости сходимости начинает проявляться лишь тогда, когда итерации уже почти сошлись и в действительности осталось сделать одну-две итерации либо вообще пора прекратить вычисления. Чтобы ощутить преимущество таких методов над методами Ньютона и секущих, нужны компьютер с очень высокой разрядностью и желание получить решение с чрезвычайно высокой точностью.

3. Вычисление корней многочленов. Вычисление корней многочленов $P_m(x)$ степени m — специальная задача, для решения которой разработано большое число эффективных алгоритмов. Укажем на один из них — *метод Мюллера*. В нем по трем ранее найденным приближениям $x^{(n-2)}$, $x^{(n-1)}$, $x^{(n)}$ строят интерполяционный многочлен второй степени $Q_2(x)$, т.е. такой многочлен, для которого $Q_2(x^{(i)}) = P_m(x^{(i)})$ при $i = n-2$, $n-1$, n . За очередное приближение $x^{(n+1)}$ принимают тот из двух корней многочлена $Q_2(x)$ (в общем случае являющихся комплексными), который расположен ближе к $x^{(n)}$. Метод Мюллера сходится с порядком $p \approx 1.839$.

Замечание. Метод Мюллера считается также эффективным методом отыскания комплексных корней функций комплексной переменной.

Вплоть до начала 60-х годов XX в. задача отыскания корней многочленов была весьма распространенной. Одна из основных причин состояла в том, что к отысканию корней характеристического многочлена сводилась одна из важнейших задач алгебры — задача вычисления собственных чисел матриц. В настоящее время существуют другие, гораздо более эффективные методы вычисления собственных значений (см. гл. 8), и, по-видимому, задача вычисления корней многочлена во многом потеряла свою актуальность.

Как было отмечено ранее, к задаче отыскания корней многочлена высокой степени с приближенно заданными коэффициентами следует отнести очень осторожно. она может оказаться плохо обусловленной (см. в гл. 3 пример Уилкинсона). Вообще, можно посоветовать избегать применения методов решения вычислительных задач, в которых используются многочлены высокой степени.

4. Гибридные алгоритмы. В последние годы получили признание алгоритмы, которые называют *гибридными* (или *регуляризованными*). Они представляют собой комбинации надежных, но медленно сходя-

шихся методов типа бисекции с недостаточно надежными, но быстро сходящимися методами типа секущих или Ньютона. Результатирующие алгоритмы обладают высокой надежностью и гарантированной сходимостью. В тех же случаях, когда в окрестности простого корня функция f — гладкая, сходимость становится сверхлинейной.

Алгоритм ZEROIN, изложенный в [106], является примером эффективного гибридного алгоритма, на каждом шаге которого принимается решение о том, какой из трех методов: бисекции, секущих или обратной квадратичной интерполяции — следует использовать для вычисления очередного приближения.

5. Δ^2 -процесс Эйткена. Пусть $\{x^{(n)}\}_{n=0}^{\infty}$ — последовательность приближений к корню \bar{x} , полученная линейно сходящимся итерационным методом (например, методом простой итерации). Использовав три последние найденные приближения $x^{(n-1)}$, $x^{(n)}$, $x^{(n+1)}$, можно получить более точное приближение к решению

$$\tilde{x}^{(n+1)} = x^{(n+1)} - \frac{(x^{(n+1)} - x^{(n)})^2}{x^{(n+1)} - 2x^{(n)} + x^{(n-1)}}.$$

Этот способ ускорения сходимости используется на завершающем этапе итераций и называется Δ^2 -процессом Эйткена.

6. Критерии окончания. Проблема выбора правильного критерия окончания итерационного процесса часто оказывается достаточно сложной. Следует также иметь в виду, что наряду с обсуждавшимся в предыдущих параграфах условием $|x^{(n)} - x^{(n-1)}| < \varepsilon_1$ нередко используется следующее условие: $|f(x^{(n)})| < \varepsilon_2$.

Глава 5

ПРЯМЫЕ МЕТОДЫ РЕШЕНИЯ СИСТЕМ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ

§ 5.1. Постановка задачи

В вычислительной линейной алгебре традиционно выделяют четыре основные задачи:

- 1) решение систем линейных алгебраических уравнений;
- 2) вычисление определителей;
- 3) нахождение обратных матриц;
- 4) определение собственных значений и собственных векторов.

В последние десятилетия к первым четырем задачам добавились еще две:

- 5) линейная задача метода наименьших квадратов;
- 6) вычисление сингулярных чисел и сингулярных векторов.

Задачи 2 и 3 обсуждаются в § 5.6, последней по порядку (но не по значению) задаче 4 посвящена гл. 8. Задачи 5 и 6 рассматриваются в § 5.12.

В основном же данная глава посвящена задаче 1, а более точно — прямым методам решения систем линейных алгебраических уравнений с вещественными коэффициентами:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1m}x_m &= b_1, \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots + a_{2m}x_m &= b_2, \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 + \dots + a_{3m}x_m &= b_3, \\ \dots & \\ a_{m1}x_1 + a_{m2}x_2 + a_{m3}x_3 + \dots + a_{mm}x_m &= b_m. \end{aligned} \tag{5.1}$$

В матричной форме записи эта система принимает вид

$$Ax = b, \tag{5.2}$$

где

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1m} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2m} \\ a_{31} & a_{32} & a_{33} & \dots & a_{3m} \\ \dots & & & & \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mm} \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_m \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_m \end{bmatrix}.$$

Итерационные методы решения системы (5.1) будут рассмотрены в гл. 6.

Уделим основное внимание задаче вычисление вектора x , являющегося решением системы (5.2), по входному данному — вектору b . Будем предполагать, что матрица A задана и является невырожденной. Известно, что в этом случае решение системы существует, единственno и устойчиво по входным данным. Это означает, что рассматриваемая задача корректна.

Хотя задача решения системы (5.1) сравнительно редко представляет самостоятельный интерес для приложений, от умения эффективно решать такие системы часто зависит сама возможность математического моделирования самых разнообразных процессов с применением компьютера. Как будет видно из дальнейшего изложения, значительная часть численных методов решения различных (в особенности — нелинейных) задач включает в себя решение систем (5.1) как элементарный шаг соответствующего алгоритма.

Пусть $x^* = (x_1^*, x_2^*, \dots, x_m^*)^T$ — приближенное решение системы (5.1). В этой и следующих главах мы будем стремиться к получению решения, для которого погрешность $e = x - x^*$ мала (количественные характеристики «величины» погрешности будут введены в следующем параграфе). Тем не менее заметим, что качество полученного решения далеко не всегда характеризуется тем, насколько мала погрешность $x - x^*$. Иногда вполне удовлетворительным является критерий малости невязки $r = b - Ax^*$. Вектор r показывает, насколько отличается правая часть системы от левой, если подставить в нее приближенное решение. Заметим, что $r = Ax - Ax^* = A(x - x^*)$ и поэтому погрешность и невязка связаны равенством

$$e = x - x^* = A^{-1}r. \quad (5.3)$$

§ 5.2. Нормы вектора и матрицы

1. Норма вектора. Решением системы линейных алгебраических уравнений является вектор $x = (x_1, x_2, \dots, x_m)^T$, который будем рассматривать как элемент векторного пространства \mathbb{R}^m . Приближенное решение $x^* = (x_1^*, x_2^*, \dots, x_m^*)^T$ и погрешность $e = x - x^* = (x_1 - x_1^*, \dots, x_m - x_m^*)^T$ также являются элементами пространства \mathbb{R}^m . Для того чтобы анализировать методы решения систем, необходимо уметь количественно оценивать «величины» векторов x^* и $x - x^*$, а также векторов b и $b - b^*$, где $b^* = (b_1^*, b_2^*, \dots, b_m^*)^T$ — вектор приближенно заданных правых частей. Удобной для этой цели количественной характеристикой является широко используемое понятие нормы вектора.

Говорят, что в \mathbb{R}^m задана норма, если каждому вектору x из \mathbb{R}^m сопоставлено вещественное число $\|x\|$, называемое *нормой вектора* x и обладающее следующими свойствами:

- 1) $\|x\| \geq 0$, причем $\|x\| = 0$ тогда и только тогда, когда $x = 0$;
- 2) $\|\alpha x\| = |\alpha| \|x\|$ для любого вектора x и любого числа α ;
- 3) $\|x + y\| \leq \|x\| + \|y\|$ для любых векторов x и y ;

последнее неравенство принято называть *неравенством треугольника*.

Заметим, что такими же свойствами обладает обычная геометрическая длина вектора в трехмерном пространстве. Свойство 3) в этом случае следует из правила сложения векторов и из того известного факта, что сумма длин двух сторон треугольника всегда больше длины третьей стороны.

Существует множество различных способов введения норм. В вычислительных методах наиболее употребительными являются следующие три нормы:

$$\|x\|_1 = \sum_{i=1}^m |x_i|, \quad \|x\|_2 = |x| = \left(\sum_{i=1}^m |x_i|^2 \right)^{1/2}, \quad \|x\|_\infty = \max_{1 \leq i \leq m} |x_i|. \quad (5.4)$$

Первые две из них являются частными случаями более общей нормы:

$$\|x\|_p = \left(\sum_{i=1}^m |x_i|^p \right)^{1/p}, \quad p \geq 1 \quad (5.5)$$

(при $p = 1$ и $p = 2$), а последняя, как можно показать, получается из нормы (5.5) предельным переходом при $p \rightarrow \infty$.

Замечание 1. Норма $\|x\|_2 = |x|$ является естественным обобщением на случай m -мерного пространства понятия длины вектора в двух- и трехмерных геометрических пространствах. Поэтому ее называют *евклидовой нормой* или *модулем* вектора x .

Замечание 2. Справедливы неравенства

$$\|x\|_\infty \leq \|x\|_2 \leq \|x\|_1 \leq m \|x\|_\infty, \quad (5.6)$$

указывающие на то, что в определенном смысле все три введенные нормы эквивалентны: каждая из них оценивается любой из двух других норм с точностью до множителя, зависящего от m .

Пример 5.1. Найдем $\|x\|_1$, $\|x\|_2$, $\|x\|_\infty$ для вектора $x = (0.12, -0.15, 0.16)^T$.

По формулам (5.4) определяем $\|x\|_1 = 0.12 + 0.15 + 0.16 = 0.43$, $\|x\|_2 = \sqrt{(0.12^2 + 0.15^2 + 0.16^2)^{1/2}} = 0.25$, $\|x\|_\infty = \max\{0.12, 0.15, 0.16\} = 0.16$. ▲

2. Скалярное произведение. Напомним, что скалярным произведением векторов $x = (x_1, x_2, \dots, x_m)^T$ и $y = (y_1, y_2, \dots, y_m)^T$ называется величина

$$(x, y) = x_1 y_1 + x_2 y_2 + \dots + x_m y_m = \sum_{i=1}^m x_i y_i. \quad (5.7)$$

Нетрудно заметить, что $\|x\|_2 = (x, x)^{1/2}$.

Когда векторы x, y имеют комплексные компоненты, скалярное произведение понимают так:

$$(x, y) = x_1 \bar{y}_1 + x_2 \bar{y}_2 + \dots + x_m \bar{y}_m = \sum_{i=1}^m x_i \bar{y}_i.$$

3. Абсолютная и относительная погрешности вектора. Далее будем всюду считать, что в пространстве m -мерных векторов R^m введена и фиксирована некоторая норма $\|x\|$ (например, одна из норм $\|x\|_p$, $1 \leq p \leq \infty$). В этом случае в качестве меры степени близости векторов x и x^* естественно использовать величину $\|x - x^*\|$, являющуюся аналогом расстояния между точками x и x^* . Введем абсолютную и относительную погрешности вектора x^* с помощью формул

$$\Delta(x^*) = \|x - x^*\|, \quad \delta(x^*) = \frac{\|x - x^*\|}{\|x\|}. \quad (5.8)$$

Выбор той или иной конкретной нормы в практических задачах диктуется тем, какие требования предъявляются к точности решения. Выбор нормы $\|x\|_1$ фактически отвечает случаю, когда малой должна быть суммарная абсолютная погрешность в компонентах решения; выбор $\|x\|_2$ соответствует критерию малости среднеквадратичной погрешности, а принятие в качестве нормы $\|x\|_\infty$ означает, что малой должна быть максимальная из абсолютных погрешностей в компонентах решения.

4. Сходимость по норме. Пусть $\{x^{(n)}\}_{n=1}^\infty$ — последовательность векторов $x^{(n)} = (x_1^{(n)}, x_2^{(n)}, \dots, x_m^{(n)})^T$. Говорят, что последовательность векторов $x^{(n)}$ сходится к вектору x при $n \rightarrow \infty$ ($x^{(n)} \rightarrow x$ при $n \rightarrow \infty$), если $\Delta(x^{(n)}) = \|x^{(n)} - x\| \rightarrow 0$ при $n \rightarrow \infty$.

Замечание. Сам факт наличия или отсутствия сходимости $x^{(n)}$ к x при $n \rightarrow \infty$ в конечномерных пространствах не зависит от выбора нормы. Известно, что из сходимости последовательности по одной из норм следует сходимость этой последовательности по любой другой норме. Например, для норм $\|x\|_1$, $\|x\|_2$, $\|x\|_\infty$ это вытекает из неравенств (5.6). Более того, $x^{(n)} \rightarrow x$ при $n \rightarrow \infty$ тогда и только тогда, когда для всех $i = 1, 2, \dots, m$ имеем $x_i^{(n)} \rightarrow x_i$ при $n \rightarrow \infty$, т.е. сходимость по норме в \mathbb{R}^m эквивалентна покомпонентной (покоординатной) сходимости.

5. Норма матрицы. Величина

$$\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \quad (5.9)$$

называется *нормой матрицы A* , подчиненной норме векторов, введенной в \mathbb{R}^m .

Заметим, что множество всех квадратных матриц размера $m \times m$ является векторным пространством. Можно показать, что введенная в этом пространстве формулой (5.9) норма обладает следующими свойствами, аналогичными свойствам нормы вектора:

- 1) $\|A\| \geq 0$, причем $\|A\| = 0$ тогда и только тогда, когда $A = 0$;
- 2) $\|\alpha A\| = |\alpha| \|A\|$ для любой матрицы A и любого числа α ;
- 3) $\|A + B\| \leq \|A\| + \|B\|$ для любых матриц A и B .

Дополнительно к этому верны следующие свойства:

- 4) $\|A \cdot B\| \leq \|A\| \cdot \|B\|$ для любых матриц A и B ;
- 5) для любой матрицы A и любого вектора x справедливо неравенство

$$\|Ax\| \leq \|A\| \cdot \|x\|. \quad (5.10)$$

Докажем, например, свойство 5). Если $\|x\| \neq 0$, то неравенство (5.10) эквивалентно неравенству $\|Ax\|/\|x\| \leq \|A\|$, справедливость которого следует из определения (5.9). Если же $\|x\| = 0$, то неравенство (5.10) превращается в верное числовое неравенство $0 \leq 0$.

Как следует из определения (5.9), каждой из векторных норм $\|x\|$ соответствует своя подчиненная норма матрицы A . Известно, в частности, что нормам $\|x\|_1$, $\|x\|_2$ и $\|x\|_\infty$ подчинены нормы $\|A\|_1$, $\|A\|_2$ и $\|A\|_\infty$, вычисляемые по формулам

$$\|A\|_1 = \max_{1 \leq j \leq m} \sum_{i=1}^m |a_{ij}|, \quad (5.11)$$

$$\|A\|_2 = \max_{1 \leq j \leq m} \sqrt{\lambda_j(A^T A)}, \quad (5.12)$$

где $\lambda_j(A^T A)$ — собственные числа¹ матрицы $A^T A$;

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^m |a_{ij}|. \quad (5.13)$$

Нормы $\|A\|_1$ и $\|A\|_\infty$ вычисляются просто (см. ниже пример 5.2). Для получения значения первой из них нужно найти сумму модулей элементов каждого из столбцов матрицы A , а затем выбрать максимальную из этих сумм. Для получения значения $\|A\|_\infty$ нужно аналогичным образом поступить со строками матрицы A .

Как правило, вычислить значение нормы $\|A\|_2$ бывает трудно, так как для этого следует искать собственные числа λ_j . Для оценки величины $\|A\|_2$ можно, например, использовать неравенство

$$\|A\|_2 \leq \|A\|_E. \quad (5.14)$$

Здесь $\|A\|_E = \sqrt{\sum_{i,j=1}^m |a_{ij}|^2}$ — величина, называемая *евклидовой нормой* матрицы A . Часто эту же величину называют *нормой Фробениуса*² и обозначают $\|A\|_F$.

Норма (5.9) имеет простую геометрическую интерпретацию. Для того чтобы ее привести, заметим, что операцию умножения матрицы A на вектор x можно рассматривать как преобразование, которое переводит вектор x в новый вектор $y = Ax$. Если значение $\|x\|$ интерпретируется как длина вектора x , то величина $\|Ax\|/\|x\|$ есть коэффициент растяжения вектора x под действием матрицы A . Таким образом, величина

$$k_{\max} = \|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \quad (5.15)$$

представляет собой максимальный коэффициент растяжения векторов под действием матрицы A . Полезно отметить, что для невырожденной матрицы A минимальный коэффициент растяжения k_{\min} отвечает норме обратной матрицы и вычисляется по формуле

$$k_{\min} = \|A^{-1}\|^{-1} = \min_{x \neq 0} \frac{\|Ax\|}{\|x\|}. \quad (5.16)$$

Заметим, что в случае $\|A\| < 1$ происходит сжатие векторов под действием матрицы A .

¹ Напомним, что число λ называется собственным числом матрицы A , если существует вектор $x \neq 0$ такой, что $Ax = \lambda x$. Каждая матрица порядка m имеет ровно m собственных чисел (вообще говоря, комплексных) с учетом их кратности.

² Фробениус Фердинанд Георг (1849—1917) — немецкий математик.

Пример 5.2. Для матрицы

$$\mathbf{A} = \begin{bmatrix} 0.1 & -0.4 & 0 \\ 0.2 & 0 & -0.3 \\ 0 & 0.1 & 0.3 \end{bmatrix}$$

найдем $\|\mathbf{A}\|_1$, $\|\mathbf{A}\|_\infty$ и оценим $\|\mathbf{A}\|_2$.

В соответствии с формулами (5.11), (5.13) и неравенством (5.14) имеем

$$\|\mathbf{A}\|_1 = \max\{0.1 + 0.2 + 0; 0.4 + 0 + 0.1; 0 + 0.3 + 0.3\} = 0.6;$$

$$\|\mathbf{A}\|_\infty = \max\{0.1 + 0.4 + 0; 0.2 + 0 + 0.3; 0 + 0.1 + 0.3\} = 0.5;$$

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_E = \left(\sum_{i,j=1}^3 a_{ij}^2 \right)^{1/2} = \sqrt{0.4} \approx 0.63. \blacksquare$$

§ 5.3. Типы используемых матриц

Эффективность вычислений в линейной алгебре существенно зависит от умения использовать специальную структуру и свойства используемых в расчетах матриц. Напомним некоторые важные типы матриц.

Квадратная матрица \mathbf{A} называется *диагональной*, если ее элементы удовлетворяют условию $a_{ij} = 0$ для $i \neq j$ (все отличные от нуля элементы расположены на главной диагонали):

$$\mathbf{A} = \begin{bmatrix} a_{11} & 0 & 0 & \dots & 0 \\ 0 & a_{22} & 0 & \dots & 0 \\ 0 & 0 & a_{33} & \dots & 0 \\ \dots & & & & \\ 0 & 0 & 0 & \dots & a_{mm} \end{bmatrix}.$$

Диагональную матрицу, у которой все элементы a_{ii} главной диагонали равны единице, называют *единичной* и обозначают буквой E :

$$E = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \dots & & & & \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}.$$

Пример 5.3. Вычислим норму единичной матрицы E .

По определению, $\|E\| = \max_{x \neq 0} \frac{\|Ex\|}{\|x\|} = \max_{x \neq 0} \frac{\|x\|}{\|x\|} = 1$. Следовательно, $\|E\| = 1$. \blacktriangle

Важную роль в численном анализе играют *треугольные матрицы*. Квадратная матрица A называется *нижней треугольной*, если все ее элементы, расположенные выше главной диагонали, равны нулю ($a_{ij} = 0$ для $i < j$). Если же равны нулю все элементы матрицы, расположенные ниже главной диагонали ($a_{ij} = 0$ для $i > j$), то она называется *верхней треугольной*.

Нижняя и верхняя треугольная матрицы имеют соответственно следующий вид:

$$\left[\begin{array}{ccccc} a_{11} & 0 & 0 & \dots & 0 \\ a_{21} & a_{22} & 0 & \dots & 0 \\ a_{31} & a_{32} & a_{33} & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mm} \end{array} \right], \quad \left[\begin{array}{ccccc} a_{11} & a_{12} & a_{13} & \dots & a_{1m} \\ 0 & a_{22} & a_{23} & \dots & a_{2m} \\ 0 & 0 & a_{33} & \dots & a_{3m} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & a_{mm} \end{array} \right].$$

Треугольные матрицы обладают рядом замечательных свойств. Например, для таких матриц определитель легко вычисляется по формуле

$$\det A = a_{11} \cdot a_{22} \cdot a_{33} \cdot \dots \cdot a_{mm}. \quad (5.17)$$

Квадратная матрица A называется *симметричной*, если она совпадает со своей транспонированной матрицей A^T (или, что то же, $a_{ij} = a_{ji}$ для всех i, j).

Будем называть симметричную матрицу A *положительно определенной* и писать $A > 0$, если для всех векторов $x \neq 0$ квадратичная форма

$$(Ax, x) = \sum_{i,j=1}^m a_{ij} x_i x_j$$

принимает положительные значения.

Обозначим через λ_{\min} и λ_{\max} минимальное и максимальное собственные значения матрицы A . Известно, что для симметричной матрицы

$$\lambda_{\min} \|x\|_2^2 \leq (Ax, x) \leq \lambda_{\max} \|x\|_2^2$$

и матрица A положительно определена тогда и только тогда, когда все ее собственные значения положительны.

Одна из трудностей практического решения систем большой размерности связана с ограниченностью оперативной памяти компьютера. Хотя объем оперативной памяти вновь создаваемых вычислительных машин растет очень быстро, тем не менее еще быстрее возрастают потребности практики в решении задач все большей размерности (напомним, что для хранения в оперативной памяти компьютера матрицы порядка m требуется m^2 машинных слов). В значительной степени ограничения на размерность решаемых систем можно снять, если использовать для хранения матрицы внешние запоминающие устройства. Однако в этом случае многократно возрастают как затраты машинного времени, так и сложность соответствующих алгоритмов. Поэтому при создании вычислительных алгоритмов линейной алгебры большое внимание уделяют способам компактного размещения элементов матриц в памяти компьютера. Заметим, что для хранения диагональной матрицы достаточно отвести массив длины m и расположить в нем элементы $a_{11}, a_{22}, \dots, a_{mm}$. Для хранения треугольной матрицы достаточно $m(m + 1)/2$ ячеек памяти, что примерно вдвое меньше места, отводимого для хранения матрицы общего вида. Столько же ячеек используется для хранения симметричной матрицы, поскольку такая матрица полностью определяется, например, заданием своей нижней треугольной части.

К счастью, приложения очень часто приводят к матрицам, в которых число ненулевых элементов много меньше общего числа элементов матрицы. Такие матрицы принято называть *разреженными*. Напротив, матрицы общего вида называют *плотными* (или *заполненными*). Разреженность матрицы является очень ценным свойством, поскольку объем информации, который следует обрабатывать и хранить в памяти компьютера, для таких матриц даже очень большого размера может оказаться не слишком большим. Для хранения всех ненулевых элементов и информации об их расположении оказывается достаточным использовать только оперативную память компьютера. Иногда элементы матрицы известны либо вычисляются по простым формулам и необходимость в их хранении отпадает.

Одним из основных источников разреженных матриц являются математические модели технических устройств, состоящих из большого числа элементов, связи между которыми локальны. Простейшие примеры таких устройств — сложные строительные конструкции и большие электрические цепи. Другой важный источник разреженности — метод конечных разностей и метод конечных элементов, используемые для решения уравнений математической физики.

Известны примеры решенных в последние годы задач, где число неизвестных достигало сотен тысяч. Естественно, это было бы невозможно, если бы соответствующие матрицы не являлись разреженными (число элементов матрицы при $m = 10^5$ равно $m \times m = 10^{10}$).

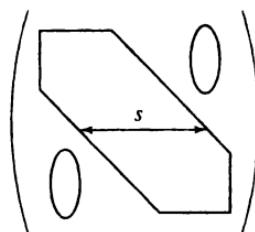


Рис. 5.1

Простой пример разреженной матрицы дает трехдиагональная матрица

$$\left[\begin{array}{ccccccc} a_{11} & a_{12} & 0 & 0 & \dots & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 & \dots & 0 & 0 \\ 0 & a_{32} & a_{33} & a_{34} & \dots & 0 & 0 \\ & & & & & & \\ 0 & 0 & 0 & 0 & \dots & a_{m-1,m-2} & a_{m-1,m-1} & a_{m-1,m} \\ 0 & 0 & 0 & 0 & \dots & 0 & a_{m,m-1} & a_{mm} \end{array} \right],$$

все ненулевые элементы которой расположены на главной и двух соседних с ней диагоналях. Число этих элементов равно $3m - 2$, что при большом m много меньше общего числа m^2 элементов матрицы.

Многие приложения приводят к системам уравнений с так называемыми ленточными матрицами. Матрица A называется ленточной с полушириной ленты, равной l , если $a_{ij} = 0$ для $|i - j| > l$. Все ненулевые элементы такой матрицы расположены на $s = 2l + 1$ ближайших к главной диагоналях матрицы; число s принято называть шириной ленты. Схематически ленточная матрица представлена на рис. 5.1. Частным случаем ленточной матрицы при $s = 3$ является трехдиагональная матрица. Ясно, что в случае $s \ll m$ ленточная матрица является разреженной.

§ 5.4. Обусловленность задачи решения системы линейных алгебраических уравнений

Оказывается, что решения различных систем линейных алгебраических уравнений обладают разной чувствительностью к погрешностям входных данных. Так же как и другие задачи (см. § 3.2), задача вычисления решения x системы уравнений

$$Ax = b \tag{5.18}$$

может быть как хорошо, так и плохо обусловленной.

Исследование обусловленности задачи начнем со случая, когда элементы матрицы A считаются заданными точно, а вектор-столбец правой части — приближенно.

Лемма 5.1. Для погрешности приближенного решения системы (5.18) справедлива оценка

$$\Delta(x^*) \leq \|A^{-1}\| \cdot \|r\|, \quad (5.19)$$

где $r = b - Ax^*$ — невязка, отвечающая x^* .

Для доказательства достаточно взять норму левой и правой частей равенства (5.3) и воспользоваться свойством (5.10). ■

Теорема 5.1. Пусть x^* — точное решение системы $Ax^* = b^*$, в которой правая часть b^* является приближением к b . Тогда верны следующие оценки абсолютной и относительной погрешностей:

$$\Delta(x^*) \leq v_\Delta \Delta(b^*), \quad (5.20)$$

$$\delta(x^*) \leq v_\delta \delta(b^*), \quad (5.21)$$

где $v_\Delta = \|A^{-1}\|$, $v_\delta(x) = \|A^{-1}\| \cdot \|b\| / \|x\| = \|A^{-1}\| \cdot \|Ax\| / \|x\|$.

Доказательство. В рассматриваемом случае $r = b - Ax^* = b - b^*$ и неравенство (5.19) принимает вид (5.20). Разделив теперь обе части неравенства (5.20) на $\|x\|$ и записав его в виде

$$\frac{\Delta(x^*)}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|b\|}{\|x\|} \frac{\Delta(b^*)}{\|b\|},$$

приходим к оценке (5.21). ■

Замечание 1. Величина $v_\Delta = \|A^{-1}\|$ для задачи (5.18) играет роль абсолютного числа обусловленности (см. § 3.2).

Замечание 2. Величина $v_\delta(x) = \|A^{-1}\| \cdot \|b\| / \|x\|$ называется естественным числом обусловленности. Она зависит от конкретного решения x и характеризует коэффициент возможного возрастания относительной погрешности этого решения, вызванного погрешностью задания правой части. Это означает, что $v_\delta(x)$ для задачи вычисления решения x системы (5.18) играет роль относительного числа обусловленности (см. § 3.2).

Замечание 3. Полученные в теореме 5.1 оценки точны в том смысле, что для системы $Ax = b$ с произвольной невырожденной матрицей A и любой заданной правой частью $b \neq 0$ найдется сколь угодно близкий к b приближенно заданный вектор $b^* \neq b$, для которого неравенства (5.20) и (5.21) превращаются в равенства.

Вычислим максимальное значение естественного числа обусловленности, используя определение (5.4) нормы матрицы:

$$\max_{x \neq 0} v_\delta(x) = \max_{x \neq 0} \frac{\|A^{-1}\| \cdot \|Ax\|}{\|x\|} = \|A^{-1}\| \cdot \|A\|. \quad (5.22)$$

Полученную величину принято называть *стандартным числом обусловленности* (или просто *числом обусловленности*) матрицы A и обозначать через $\text{cond}(A)$ или $\text{cond}(A)$. Таким образом,

$$v(A) = \text{cond}(A) = \|A^{-1}\| \cdot \|A\|. \quad (5.23)$$

Сформулируем важное следствие из теоремы 5.1.

Следствие. В условиях теоремы 5.1 справедлива оценка

$$\delta(x^*) \leq \text{cond}(A) \cdot \delta(b^*). \quad (5.24)$$

Для ее доказательства достаточно воспользоваться оценкой (5.21) и заметить, что в силу равенства (5.22) верно неравенство $v_\delta \leq \text{cond}(A)$. ■

Замечание. Оценка (5.24) точна в том смысле, что для системы (5.18) с произвольной невырожденной матрицей A найдутся правая часть $b \neq 0$ (и отвечающее этой правой части решение x) и сколь угодно близкий к b приближенно заданный вектор $b^* \neq b$ такие, что неравенство (5.24) превращается в равенство.

Величина $\text{cond}(A)$ является широко используемой количественной мерой обусловленности системы $Ax = b$. В частности, систему и матрицу A принято называть *плохо обусловленными*, если $\text{cond}(A) \gg 1$. В силу оценки (5.24) и последнего замечания для такой системы существуют решения, обладающие чрезвычайно высокой чувствительностью к малым погрешностям задания входного данного b . Тем не менее заметим, что не для всякого решения x коэффициент $v_\delta(x)$ роста относительной погрешности достигает значений, близких к максимально возможному значению $\text{cond}(A)$.

Отметим следующие свойства числа обусловленности.

1°. Для единичной матрицы $\text{cond}(E) = 1$.

Доказательство. Пользуясь тем, что $E^{-1} = E$ и $\|E\| = 1$ (см. пример 5.3), получаем $\text{cond}(E) = \|E^{-1}\| \cdot \|E\| = 1$. ■

2°. Справедливо неравенство $\text{cond}(A) \geq 1$.

Доказательство. Из равенства $E = A \cdot A^{-1}$, свойства 4° норм матриц и равенства $\|E\| = 1$ следует, что $1 = \|E\| \leq \|A^{-1}\| \cdot \|A\| = \text{cond}(A)$. ■

3°. Число обусловленности матрицы A не меняется при умножении матрицы на произвольное число $a \neq 0$.

Доказательство. Заметим, что $(\alpha A)^{-1} = \alpha^{-1} A^{-1}$. Поэтому
 $\text{cond}(\alpha A) = \|\alpha A\| \cdot \|(\alpha A)^{-1}\| = |\alpha| \cdot \|A\| |\alpha|^{-1} \|A^{-1}\| = \text{cond}(A)$. ■

Замечание. Пользуясь приведенной в § 5.2 геометрической интерпретацией норм матриц A и A^{-1} (см. формулы (5.15) и (5.16)), число обусловленности можно интерпретировать как отношение максимального коэффициента растяжения векторов под действием матрицы A к минимальному коэффициенту $\text{cond}(A) = k_{\max}/k_{\min}$.

Величина $\text{cond}(A)$ зависит, вообще говоря, от выбора нормы векторов в пространстве \mathbb{R}^m . Фактически это есть зависимость максимального коэффициента роста погрешности от способа измерения величины входных данных и решения. В частности, выбору нормы $\|x\|_p$ ($1 \leq p \leq \infty$) отвечает $\text{cond}_p(A) = \|A^{-1}\|_p \cdot \|A\|_p$.

Пример 5.4. Вычислим $\text{cond}_{\infty}(A)$ для матрицы

$$A = \begin{bmatrix} 1.03 & 0.991 \\ 0.991 & 0.943 \end{bmatrix}. \quad (5.25)$$

Сначала найдем обратную матрицу

$$A^{-1} \approx \begin{bmatrix} -87.4 & 91.8 \\ 91.8 & -95.4 \end{bmatrix}.$$

Тогда $\text{cond}_{\infty}(A) = \|A\|_{\infty} \cdot \|A^{-1}\|_{\infty} \approx 2.021 \cdot 187.2 \approx 378$. Если входные данные для системы уравнений с матрицей (5.25) содержат относительную погрешность порядка 0.1—1 %, то систему можно расценить как плохо обусловленную. ▲

Пример 5.5. Рассмотрим систему уравнений

$$\begin{aligned} 1.03x_1 + 0.991x_2 &= 2.51, \\ 0.991x_1 + 0.943x_2 &= 2.41 \end{aligned} \quad (5.26)$$

с матрицей (5.25). Ее решением является $x_1 \approx 1.981$, $x_2 \approx 0.4735$. Правая часть системы известна в лучшем случае с точностью до 0.005, если считать, что числа 2.51 и 2.41 получены округлением «истинных» значений при вводе в память трехзначного десятичного компьютера. Как влияет погрешность во входных данных такого уровня на погрешность решения? Возьмутим каждую из компонент вектора правой части $b = (2.51, 2.41)^T$ на 0.005, взяв $b^* = (2.505, 2.415)^T$. Решением системы,

отвечающим \mathbf{b}^* , является теперь $x_1^* \approx 2.877$, $x_2^* \approx -0.4629$. Таким образом, решение оказалось полностью искаженным. Относительная погрешность задания правой части $\delta(\mathbf{b}^*) = \|\mathbf{b} - \mathbf{b}^*\|_\infty / \|\mathbf{b}\|_\infty = 0.005/2.51 \approx \approx 0.2\%$ привела к относительной погрешности решения $\delta(\mathbf{x}^*) = \|\mathbf{x} - \mathbf{x}^*\|_\infty / \|\mathbf{x}\|_\infty \approx 0.9364/1.981 \approx 47.3\%$. Следовательно, погрешность возросла примерно в 237 раз.

Можно ли внести в правую часть системы (5.26) такую погрешность, чтобы получить существенно большее, чем 237, значение коэффициента роста погрешности? Вычислим естественное число обусловленности, являющееся максимальным значением рассматриваемого коэффициента, отвечающим решению $\mathbf{x} \approx (1.981, 0.4735)^T$ и получим $v_\delta(\mathbf{x}) = \|A^{-1}\|_\infty \|\mathbf{b}\|_\infty / \|\mathbf{x}\|_\infty \approx 187.2 \cdot 2.51 / 1.981 \approx 237$. Таким образом, на поставленный вопрос следует ответить отрицательно.

Можно дать следующую геометрическую интерпретацию рассмотренного примера. Каждому уравнению системы (5.26) соответствует прямая на плоскости Ox_1x_2 . По коэффициентам при x_1 и x_2 в этих уравнениях видно, что прямые почти параллельны. Так как вблизи точки пересечения прямые сливаются, то даже незначительная погрешность в задании положения этих прямых существенно меняет положение точки пересечения (рис. 5.2). ▲

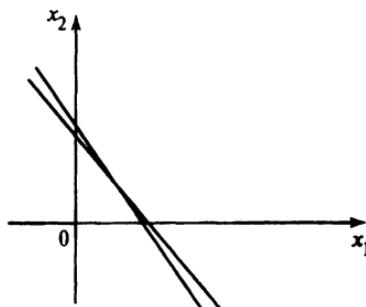


Рис. 5.2

Пример 5.6. Традиционным примером очень плохо обусловленной матрицы является матрица Гильберта¹ — матрица \mathbf{H} с элементами $h_{ij} = 1/(i + j - 1)$.

¹ Давид Гильберт (1862—1943) — немецкий математик, исследования которого оказали большое влияние на развитие современной математики.

Таблица 5.1

Порядок матрицы Гильберта	2	3	4	5	6	7	8	9	10
Приближенное значение числа обусловленности	$2 \cdot 10^1$	$5 \cdot 10^2$	$2 \cdot 10^4$	$5 \cdot 10^5$	$2 \cdot 10^7$	$5 \cdot 10^8$	$2 \cdot 10^{10}$	$5 \cdot 10^{11}$	$2 \cdot 10^{13}$

Из табл. 5.1, заимствованной из [107], видно, что для матрицы H даже сравнительно невысокого порядка число обусловленности оказывается чрезвычайно большим. ▲

До сих пор мы предполагали, что матрица A задана точно. Однако на практике это часто не так. Как выясняется, введенная выше величина $\text{cond}(A)$ характеризует также и чувствительность решений системы к малым погрешностям задания элементов матрицы A . В подтверждение сказанного приведем следующий результат.

Теорема 5.2. Пусть x^* — точное решение системы $A_*x^* = b$ с приближенно заданной матрицей A_* . Тогда верна следующая оценка относительной погрешности:

$$\delta^*(x^*) \leq \text{cond}(A) \cdot \delta(A_*), \quad (5.27)$$

где $\delta^*(x^*) = \|x - x^*\| / \|x^*\|$, $\delta(A_*) = \|A - A_*\| / \|A\|$.

Доказательство. В данном случае невязка r имеет вид $r = b - Ax^* = A_*x^* - Ax^* = (A_* - A)x^*$. Применяя неравенство (5.19), получаем цепочку неравенств

$$\begin{aligned} \delta^*(x^*) &= \frac{\|x - x^*\|}{\|x^*\|} \leq \frac{\|A^{-1}\| \cdot \|(A_* - A)x^*\|}{\|x^*\|} \leq \frac{\|A^{-1}\| \cdot \|A_* - A\| \cdot \|x^*\|}{\|x^*\|} = \\ &= \text{cond}(A) \cdot \delta(A_*). \blacksquare \end{aligned}$$

Следствие. В условиях теоремы 5.2 справедливо приближенное неравенство

$$\delta(x^*) \lesssim \text{cond}(A) \cdot \delta(A_*). \quad (5.28)$$

Замечание 1. Когда с погрешностью заданы как правая часть системы, так и матрица (т.е. x^* является решением системы $A_*x^* = b^*$), причем $\text{cond}(A) \cdot \delta(A_*) \ll 1$, можно доказать справедливость неравенства

$$\delta(x^*) \lesssim \text{cond}(A)(\delta(b^*) + \delta(A_*)).$$

Замечание 2. Распространенным является представление о том, что по значению определителя матрицы A можно судить о степени близости системы уравнений к вырожденной или об обусловленности системы. Для того чтобы убедиться в ошибочности этого мнения, умножим каждое из уравнений системы (5.1) на постоянную $\alpha \neq 0$. Ясно, что такое преобразование никак не меняет решение системы и его чувствительность к малым относительным погрешностям в данных. Однако определитель умножается на число α^n и поэтому с помощью выбора α может быть сделан как угодно большим или малым. Подчеркнем, что число обусловленности $\text{cond}(A)$ при таком преобразовании системы не меняется в силу свойства 3°.

Замечание 3. Вычисление чисел обусловленности $v_\delta(x) = \|A^{-1}\| \cdot \|b\| / \|x\|$ и $\text{cond}(A) = \|A^{-1}\| \cdot \|A\|$ непосредственно по указанным формулам предполагает предварительное вычисление обратной матрицы A^{-1} . Вследствие большой трудоемкости этой операции (как показано в § 5.6, для ее выполнения в общем случае требуется примерно $2m^3$ арифметических операций) на практике избегают такого способа вычисления. При этом важно отметить, что в большинстве случаев достаточно лишь грубой оценки числа обусловленности с точностью до порядка. С эффективными методами, дающими оценки величин $v_\delta(x)$ и $\text{cond}(A)$, можно познакомиться в [81, 106, 54].

Проверить чувствительность решения системы $Ax = b$ к погрешностям можно и экспериментально. Для этого достаточно решить задачу несколько раз с близкими к b правыми частями $b^{(1)}, b^{(2)}, \dots, b^{(n)}$. Можно

ожидать, что величина $\tilde{v}_\delta = \max_{1 \leq l \leq n} \frac{\delta(x^{(l)})}{\delta(b^{(l)})}$ даст оценку значения $v_\delta(x)$.

Во всяком случае она дает оценку снизу, так как $\tilde{v}_\delta \leq v_\delta(x) \leq \text{cond}(A)$.

§ 5.5. Метод Гаусса

Рассмотрим один из самых распространенных методов решения систем линейных алгебраических уравнений — *метод Гаусса*¹. Этот метод (который называют также *методом последовательного исключения неизвестных*) известен в различных вариантах уже более 2000 лет.

¹ Карл Фридрих Гаусс (1777—1855) — немецкий математик и физик, работы которого оказали большое влияние на дальнейшее развитие высшей алгебры, геометрии, теории чисел, теории электричества и магнетизма.

Во всяком случае, его использовали китайские математики примерно за 250 лет до новой эры.

Вычисления с помощью метода Гаусса состоят из двух основных этапов, называемых *прямым ходом* и *обратным ходом* (*обратной подстановкой*). Прямой ход метода Гаусса заключается в последовательном исключении неизвестных из системы (5.1) для преобразования ее к эквивалентной системе с верхней треугольной матрицей. Вычисления значений неизвестных производят на этапе обратного хода.

1. Схема единственного деления. Рассмотрим сначала простейший вариант метода Гаусса, называемый **схемой единственного деления**.

Прямой ход состоит из $m - 1$ шагов исключения.

1-й шаг. Целью этого шага является исключение неизвестного x_1 из уравнений с номерами $i = 2, 3, \dots, m$. Предположим, что коэффициент $a_{11} \neq 0$. Будем называть его **главным (или ведущим) элементом 1-го шага**.

Найдем величины

$$\mu_{ij} = a_{ij}/a_{11} \quad (i = 2, 3, \dots, m), \quad (5.29)$$

называемые множителями 1-го шага. Вычтем последовательно из второго, третьего, ..., m -го уравнений системы (5.1) первое уравнение, умноженное соответственно на $\mu_{21}, \mu_{31}, \dots, \mu_{m1}$. Это позволит обратить в нуль коэффициенты при x_1 во всех уравнениях, кроме первого. В результате получим эквивалентную систему

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1m}x_m &= b_1, \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2m}^{(1)}x_m &= b_2^{(1)}, \\ a_{32}^{(1)}x_2 + a_{33}^{(1)}x_3 + \dots + a_{3m}^{(1)}x_m &= b_3^{(1)}, \end{aligned} \quad (5.30)$$

$$a_{m2}^{(1)}x_2 + a_{m3}^{(1)}x_3 + \dots + a_{mm}^{(1)}x_m = b_m^{(1)},$$

в которой $a_{ij}^{(1)}$ и $b_i^{(1)}$ ($i, j = 2, 3, \dots, m$) вычисляются по формулам

$$a_{ii}^{(1)} = a_{ii} - \mu_{ii} a_{1i}, \quad b_i^{(1)} = b_i - \mu_{ii} b_1. \quad (5.31)$$

2-й шаг. Целью этого шага является исключение неизвестного x_2 из уравнений с номерами $i = 3, 4, \dots, m$. Пусть $a_{22}^{(1)} \neq 0$, где $a_{22}^{(1)}$ — коэффициент, называемый **главным** (или **ведущим**) элементом 2-го шага. Вычислим множители 2-го шага

$$\mu_{i2} = a_{i2}^{(1)}/a_{22}^{(1)} \quad (i = 3, 4, \dots, m)$$

и вычтем последовательно из третьего, четвертого, ..., m -го уравнений системы (5.30) второе уравнение, умноженное соответственно на $\mu_{32}, \mu_{42}, \dots, \mu_{m2}$. В результате получим систему

Здесь $a_{ij}^{(2)}$ и $b_i^{(2)}$ ($i, j = 3, 4, \dots, m$) вычисляются по формулам

$$a_{ii}^{(2)} = a_{ii}^{(1)} - \mu_{i2} a_{2i}^{(1)}, \quad b_i^{(2)} = b_i^{(1)} - \mu_{i2} b_2^{(1)}.$$

Аналогично проводятся остальные шаги. Опишем очередной k -й шаг. k -й шаг. В предположении, что главный (ведущий) элемент k -го шага $a_{kk}^{(k-1)}$ отличен от нуля, вычислим множители k -го шага

$$\mu_{ik} = a_{ik}^{(k-1)} / a_{kk}^{(k-1)} \quad (i = k+1, \dots, m)$$

и вычтем последовательно из $(k+1)$ -го, ..., m -го уравнений полученной на предыдущем шаге системы k -е уравнение, умноженное соответственно на $\mu_{k+1,k}, \mu_{k+2,k}, \dots, \mu_{mk}$.

После $(m - 1)$ -го шага исключения получим систему уравнений

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m &= b_1, \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2m}^{(1)}x_m &= b_2^{(1)}, \\ a_{33}^{(2)}x_3 + \dots + a_{3m}^{(2)}x_m &= b_3^{(2)}, \\ \dots &\dots \\ a_{mm}^{(m-1)}x_m &= b_m^{(m-1)}, \end{aligned} \quad (5.33)$$

матрица $A^{(m-1)}$ которой является верхней треугольной. На этом вычисления прямого хода заканчиваются.

Обратный ход. Из последнего уравнения системы (5.33) находим x_m . Подставляя найденное значение x_m в предпоследнее уравнение, получаем x_{m-1} . Осуществляя обратную подстановку, далее последова-

тельно находим $x_{m-2}, x_{m-3}, \dots, x_1$. Вычисления неизвестных здесь проводятся по формулам

$$\begin{aligned} x_m &= b_m^{(m-1)} / a_{mm}^{(m-1)}, \\ x_k &= \frac{b_k^{(k-1)} - a_{k,k+1}^{(k-1)}x_{k+1} - \dots - a_{km}^{(k-1)}x_m}{a_{kk}^{(k-1)}} \quad (k = m-1, \dots, 1). \end{aligned} \quad (5.34)$$

Трудоемкость метода. Оценим число арифметических операций, необходимых для реализации схемы единственного деления.

Вычисления 1-го шага исключения по формулам (5.29), (5.31) требуют выполнения $m-1$ деления, $(m-1)m$ умножений и $(m-1)m$ вычитаний, т.е. общее число арифметических операций составляет $Q_1 = 2(m-1)^2 + 3(m-1)$. Аналогично, на 2-м шаге требуется $Q_2 = 2(m-2)^2 + 3(m-2)$ операций, а на k -м шаге — $Q_k = 2(m-k)^2 + 3(m-k)$ операций.

Подсчитаем теперь приближенно общее число Q арифметических операций прямого хода, считая размерность системы m достаточно большой:

$$\begin{aligned} Q &= \sum_{k=1}^{m-1} Q_k = 2 \sum_{k=1}^{m-1} (m-k)^2 + 3 \sum_{k=1}^{m-1} (m-k) = 2 \sum_{k=1}^{m-1} k^2 + 3 \sum_{k=1}^{m-1} k = \\ &= \frac{2(m-1)m(2m-1)}{6} + \frac{3(m-1)m}{2} \approx \frac{2}{3} m^3. \end{aligned}$$

Как нетрудно видеть, для реализации обратного хода по формулам (5.34) нужно всего m^2 операций, что при больших m пренебрежимо мало по сравнению с числом операции прямого хода.

Таким образом, для реализации метода Гаусса требуется примерно $(2/3)m^3$ арифметических операций, причем подавляющее число этих действий совершается на этапе прямого хода.

Пример 5.7. Методом Гаусса решим систему

$$\begin{aligned} 10x_1 + 6x_2 + 2x_3 &= 25, \\ 5x_1 + x_2 - 2x_3 + 4x_4 &= 14, \\ 3x_1 + 5x_2 + x_3 - x_4 &= 10, \\ 6x_2 - 2x_3 + 2x_4 &= 8. \end{aligned} \quad (5.35)$$

Прямой ход.

1-й шаг. Вычислим множители $\mu_{21} = a_{21}/a_{11} = 5/10 = 0.5$, $\mu_{31} = a_{31}/a_{11} = 3/10 = 0.3$, $\mu_{41} = a_{41}/a_{11} = 0/10 = 0$. Вычитая из второго,

третьего и четвертого уравнений системы (5.35) первое уравнение, умноженное на μ_{21} , μ_{31} и μ_{41} соответственно получаем:

$$\begin{aligned} 10x_1 + 6x_2 + 2x_3 &= 25, \\ -2x_2 - 3x_3 + 4x_4 &= 1.5, \\ 3.2x_2 + 0.4x_3 - x_4 &= 2.5, \\ 6x_2 - 2x_3 + 2x_4 &= 8. \end{aligned} \quad (5.36)$$

2-й шаг. Вычислим множители $\mu_{32} = a_{32}^{(1)} / a_{22}^{(1)} = 3.2 / (-2) = -1.6$, $\mu_{42} = 6 / (-2) = -3$. Вычитая из третьего и четвертого уравнений системы (5.36) второе уравнение, умноженное на μ_{32} и μ_{42} соответственно, приходим к системе

$$\begin{aligned} 10x_1 + 6x_2 + 2x_3 &= 25, \\ -2x_2 - 3x_3 + 4x_4 &= 1.5, \\ -4.4x_3 + 5.4x_4 &= 4.9, \\ -11x_3 + 14x_4 &= 12.5. \end{aligned} \quad (5.37)$$

3-й шаг. Вычисляя множитель $\mu_{43} = (-11) / (-4.4) = 2.5$ и вычитая из четвертого уравнения системы (5.37) третье уравнение, умноженное на μ_{43} , приводим систему к треугольному виду:

$$\begin{aligned} 10x_1 + 6x_2 + 2x_3 &= 25, \\ -2x_2 - 3x_3 + 4x_4 &= 1.5, \\ -4.4x_3 + 5.4x_4 &= 4.9, \\ 0.5x_4 &= 0.25. \end{aligned} \quad (5.38)$$

Обратный ход. Из последнего уравнения системы находим $x_4 = 0.5$. Подставляя значение x_4 в третье уравнение, находим $x_3 = (4.9 - 5.4x_4) / (-4.4) = -0.5$. Продолжая далее обратную подстановку, получаем $x_2 = (1.5 + 3x_3 - 4x_4) / (-2) = 1$, $x_1 = (25 - 6x_2 - 2x_3) / 10 = 2$. Итак, $x_1 = 2$, $x_2 = 1$, $x_3 = -0.5$, $x_4 = 0.5$.

Результаты вычислений приведены в табл. 5.2. ▲

Необходимость выбора главных элементов. Заметим, что вычисление множителей, а также обратная подстановка предполагают деление на главные элементы $a_{kk}^{(k-1)}$. Поэтому если один из главных

Таблица 5.2

Наименование	a_{11}	a_{12}	a_{13}	a_{14}	b_i	μ_{ij}, x_i
Исходная система	10	6	2	0	25	—
	5	1	-2	4	14	—
	3	5	1	-1	10	—
	0	6	-2	2	8	—
1-й шаг прямого хода	10	6	2	0	25	—
	0	-2	-3	4	1.5	0.5
	0	3.2	0.4	-1	2.5	0.3
	0	6	-2	2	8	0
2-й шаг прямого хода	10	6	2	0	25	—
	0	-2	-3	4	1.5	—
	0	0	-4.4	5.4	4.9	-1.6
	0	0	-11	14	7.5	-3
3-й шаг прямого хода и обратный ход	10	6	2	0	25	2
	0	-2	-3	4	1.5	1
	0	0	-4.4	5.4	4.9	-0.5
	0	0	0	0.5	0.25	0.5

элементов оказывается равным нулю, то схема единственного деления не может быть реализована. Здравый смысл подсказывает, что и в ситуации, когда все главные элементы отличны от нуля, но среди них есть близкие к нулю, возможен неконтролируемый рост погрешности.

Пример 5.8. Используя метод Гаусса, решим систему уравнений

$$\begin{aligned} 2x_1 - 9x_2 + 5x_3 &= -4, \\ 1.2x_1 - 5.3999x_2 + 6x_3 &= 0.6001, \\ x_1 - x_2 - 7.5x_3 &= -8.5 \end{aligned} \quad (5.39)$$

на 6-разрядном десятичном компьютере.

Прямой ход.

1-й шаг. Вычисляем множители $\mu_{21} = 0.6$ и $\mu_{31} = 0.5$ и преобразуем систему к виду

$$\begin{aligned} 2x_1 - 9x_2 + 5x_3 &= -4, \\ 0.0001x_2 + 3x_3 &= 3.0001, \\ 3.5x_2 - 10x_3 &= -6.5. \end{aligned} \quad (5.40)$$

Все вычисления на этом шаге выполняются без округлений.

2-й шаг. После вычисления множителя $\mu_{32} = 3.5/0.0001 = 35\ 000$ последнее уравнение системы должно быть преобразовано к виду

$a_{33}^{(2)}x_3 = b_3^{(2)}$, где $a_{33}^{(2)} = -10 - 3 \cdot \mu_{32} = -105\ 010$, $b_3^{(2)} = -6.5 - 3.0001 \cdot \mu_{32} = -105\ 010$. Однако на используемом компьютере будет получено уравнение

$$-105\ 010x_3 = -105\ 011. \quad (5.41)$$

Действительно, коэффициент $a_{33}^{(2)}$ определяется точно, так как при его вычислении не возникает чисел, мантиссы которых имеют более шести разрядов. В то же время при вычислении $b_3^{(2)}$ умножение коэффициента 3.0001 на μ_{32} дает 7-разрядное число 105 003.5, после округления которого до шести разрядов получится 105 004. Вычисление $b_3^{(2)}$ завершается выполнением операции вычитания: $b_3^{(2)} \approx -6.5 - 105\ 004 = -105\ 010.5$. После округления последнего числа до шести разрядов мантиссы приходим к уравнению (5.41).

Обратный ход. Из уравнения (5.41) находим $x_3 \approx 1.00001$. Сравнение с истинным значением $x_3 = 1$ показывает, что эта величина получена с очень высокой для используемого компьютера точностью.

Дальнейшие вычисления дают

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

$$x_1 = (-4 + 9x_2 - 5x_3)/2 = (-4 + 6.3 - 5.00005)/2 = -1.350025.$$

После округления имеем $x_1 = -1.35003$.

Как нетрудно видеть, найденные значения неизвестных имеют мало общего с истинными значениями решения $x_1 = 0$, $x_2 = 1$, $x_3 = 1$.

В чем же причина появления такой значительной погрешности? Говорить о накоплении ошибок округления не приходится, так как всего было выполнено 28 арифметических операций и лишь в четырех случаях потребовалось округление. Предположение о плохой обусловленности системы не подтверждается; вычисление $\text{cond}_{\infty}(A)$ дает значение, равное примерно 100.

В действительности причина состоит в использовании на 2-м шаге малого ведущего элемента $a_{22}^{(2)} = 0.0001$. Следствием этого стало появление большого множителя μ_{32} и существенное возрастание коэффициентов в последнем уравнении системы. ▲

Таким образом, изложенный выше вариант метода Гаусса (схема единственного деления) оказался некорректным и, следовательно, непригодным для вычислений на компьютере. Этот метод может привести

к делению на нуль (если $a_{kk}^{(k-1)} = 0$ при некотором k), и вычисления по нему могут оказаться неустойчивыми.

2. Метод Гаусса с выбором главного элемента по столбцу (схема частичного выбора).

Описание метода. На k -м шаге прямого хода коэффициенты уравнений системы с номерами $i = k + 1, \dots, m$ преобразуются по формулам

$$a_{ij}^{(k)} = a_{ij}^{(k-1)} - \mu_{ik} a_{kj}^{(k-1)}, \quad b_i^{(k)} = b_i^{(k-1)} - \mu_{ik} b_k^{(k-1)}, \quad i = k + 1, \dots, m. \quad (5.42)$$

Интуитивно ясно, что во избежание сильного роста коэффициентов системы и связанных с этим погрешностей нельзя допускать появления больших множителей μ_{ik} .

В методе Гаусса с выбором главного элемента¹ по столбцу гарантируется, что $|\mu_{ik}| \leq 1$ для всех $k = 1, 2, \dots, m-1$ и $i = k+1, \dots, m$. Отличие этого варианта метода Гаусса от схемы единственного деления заключается в том, что на k -м шаге исключения в качестве главного элемента выбирают максимальный по модулю коэффициент $a_{i_k k}^{(k-1)}$ при неизвестной x_k в уравнениях с номерами $i = k, k+1, \dots, m$. Затем соответствующее выбранному коэффициенту уравнение с номером i_k меняют местами с k -м уравнением системы для того, чтобы главный элемент занял место коэффициента $a_{kk}^{(k-1)}$.

После этой перестановки исключение неизвестного x_k производят, как в схеме единственного деления.

Пример 5.9. Решим систему уравнений (5.39) методом Гаусса с выбором главного элемента по столбцу на 6-разрядном десятичном компьютере.

Прямой ход.

1-й шаг. Максимальный в первом столбце элемент матрицы находится в первой строке, поэтому перестановка уравнений не нужна. Здесь 1-й шаг проводится точно так же, как и в примере 5.8.

2-й шаг. Среди элементов $a_{22}^{(1)} = 0.0001$ и $a_{32}^{(1)} = 3.5$ матрицы системы (5.40) максимальный по модулю принадлежит третьему уравнению. Меняя местами второе и третье уравнения, получаем систему

$$2x_1 - 9x_2 + 5x_3 = -4,$$

$$3.5x_2 - 10x_3 = -6.5,$$

$$0.0001x_2 + 3x_3 = 3.0001.$$

¹ Выбор главного элемента иногда называют *пivотированием*.

После вычисления $\mu_{32} = 0.0001/3.5 \approx 2.85714 \cdot 10^{-5}$ последнее уравнение системы преобразуется к виду $3.00029x_3 = 3.00029$.

Обратный ход. Из последнего уравнения находим $x_3 = 1$. Далее, имеем $x_2 = (-6.5 + 10x_3)/3.5 = 1$, $x_1 = (-4 + 9x_2 - 5x_3)/2 = (-4 + 9 - 5)/2 = 0$. В данном случае ответ получился точным. \blacktriangle

Заметим, что дополнительная работа по выбору главных элементов в схеме частичного выбора требует порядка m^2 действий, что практически не влияет на общую трудоемкость метода.

Вычислительная устойчивость схемы частичного выбора. Детальное исследование метода Гаусса показывает, что действительной причиной неустойчивости схемы единственного деления является возможность неограниченного роста элементов промежуточных матриц $A^{(1)}, A^{(2)}, \dots, A^{(m-1)}$ в процессе прямого хода. Так как на k -м шаге схемы частичного выбора $|\mu_{kk}| \leq 1$, то для вычисленных по формулам (5.42) элементов $a_{ij}^{(k)}$ справедлива оценка $|a_{ij}^{(k)}| \leq |a_{ii}^{(k-1)}| + |a_{kj}^{(k-1)}|$. Следовательно, максимальное по модулю значение элементов матрицы возрастает на одном шаге не более чем в 2 раза и в самом неблагоприятном случае $m-1$ шаг прямого хода даст коэффициент роста $\varphi(m) = 2^{m-1}$.

Гарантия ограниченности роста элементов матрицы делает схему частичного выбора вычислительно устойчивой. Более того, для нее оказывается справедливой следующая оценка погрешности:

$$\delta(x^*) \leq f(m) \operatorname{cond}_E(A) \cdot \varepsilon_m. \quad (5.43)$$

Здесь x^* — вычисленное на компьютере решение системы; $\delta(x^*) = \|x - x^*\|_2 / \|x\|_2$ — его относительная погрешность; $\operatorname{cond}_E(A) = \|A\|_E \|A^{-1}\|_E$ — число обусловленности матрицы A ; ε_m — машинное эпсилон; наконец, $f(m) = C(m)\varphi(m)$, где $C(m)$ — некоторая медленно растущая функция, зависящая от порядка m системы (типа степенной функции с небольшим показателем).

Наличие в оценке (5.43) множителя $\varphi(m) = 2^{m-1}$ указывает на то, что при большом m схема частичного выбора может оказаться плохо обусловленной и возможна существенная потеря точности. Однако практика матричных вычислений показывает, что существенный рост элементов матрицы происходит крайне редко. В подавляющем большинстве случаев действительное значение коэффициента роста не превышает 8—10. Если система хорошо обусловлена, то погрешность вычисленного решения оказывается, как правило, малой.

Иногда для проверки качества приближенного решения x^* вычисляют невязку $r = b - Ax^*$ и о степени близости приближенного решения

к точному пытаются судить по тому, насколько мала невязка. Этот метод недостаточно надежен по отношению к схеме частичного выбора, так как известно, что она гарантированно дает малые невязки. Более точно это утверждение можно сформулировать так: справедлива оценка

$$\|\mathbf{r}\|_2 \leq f(m) \|\mathbf{A}\|_E \|\mathbf{x}\|_2 \varepsilon_m, \quad (5.44)$$

где $f(m)$ то же, что и в оценке (5.43). Заметим, что в неравенство (5.44) не входит число обусловленности.

3. Метод Гаусса с выбором главного элемента по всей матрице (схема полного выбора). В этой схеме допускается нарушение естественного порядка исключения неизвестных.

На 1-м шаге метода среди элементов a_{ij} определяют максимальный по модулю элемент $a_{i_1 j_1}$. Первое уравнение системы и уравнение с номером i_1 меняют местами. Далее стандартным образом производят исключение неизвестного x_{j_1} из всех уравнений, кроме первого.

На k -м шаге метода среди коэффициентов $a_{ij}^{(k-1)}$ при неизвестных в уравнениях системы с номерами $i = k, \dots, m$ выбирают максимальный по модулю коэффициент $a_{i_k j_k}^{(k-1)}$. Затем k -е уравнение и уравнение, содержащее найденный коэффициент, меняют местами и исключают неизвестное x_{j_k} из уравнений с номерами $i = k + 1, \dots, m$.

На этапе обратного хода неизвестные вычисляют в следующем порядке: $x_{j_m}, x_{j_{m-1}}, \dots, x_{j_1}$.

Пример 5.10. Решим систему (5.39), используя схему полного выбора на 6-разрядном десятичном компьютере.

Прямой ход.

1-й шаг. Максимальный по модулю элемент $a_{12} = -9$ содержится в первом уравнении, поэтому перестановка уравнений не нужна. Исключаем неизвестное x_2 из второго и третьего уравнений, используя множители $\mu_{21} = -5.3999/(-9) \approx 0.599989$ и $\mu_{31} = -1/(-9) \approx 0.111111$. В результате получаем систему

$$\begin{aligned} 2x_1 - 9x_2 + 5x_3 &= -4, \\ 0.000022x_1 + 3.00006x_3 &= 3.00006, \\ 0.777778x_1 - 8.05556x_3 &= -8.05556. \end{aligned}$$

2-й шаг. Среди коэффициентов при неизвестных во втором и третьем уравнениях максимальным по модулю является коэффициент $a_{33}^{(1)} = -8.05556$. Переставляя местами второе и третье уравнения и исключая неизвестное x_3 (соответствующий множитель $\mu_{32} = 3.00006/(-8.05556) \approx -0.372421$), приходим к системе

$$\begin{aligned} 2x_1 - 9x_2 + 5x_3 &= -4, \\ 0.777778x_1 - 8.05556x_3 &= -8.05556, \\ 0.289683x_1 &= 2.89240 \cdot 10^{-7}. \end{aligned}$$

Обратный ход. Из последнего уравнения находим $x_1 = 2.89240 \cdot 10^{-7}/0.289683 \approx 9.98470 \cdot 10^{-7}$. Далее, имеем $x_3 = (-8.05556 - 0.777778x_1)/(-8.05556) \approx 1.00000$, $x_2 = (-4 - 2x_1 - 5x_3)/(-9) \approx 1.00000$.

Округляя найденные значения до пяти цифр после десятичной точки, получим ответ: $x_1 = 0.00000$, $x_2 = 1.00000$, $x_3 = 1.00000$. Заметим, что в данном случае получено решение, совпадающее с точным. ▲

Схема полного выбора по сравнению со схемой частичного выбора дает существенное замедление роста элементов матрицы. Доказано, что для нее коэффициент роста $\phi(m)$, входящий в оценку (5.43), не превышает значения $m^{1/2} (2^1 \cdot 3^{1/2} \cdot 4^{1/3} \dots m^{1/(m-1)})^{1/2} \leq 1.8m^{1/2}m^{0.25\ln m}$ (что значительно меньше соответствующего значения $\phi(m) = 2^{m-1}$ для схемы частичного выбора). Более того, длительное время существовала гипотеза Уилкинсона, согласно которой для схемы полного выбора $\phi(m) \leq m$. Только в 1991 г. была найдена¹ матрица 13-го порядка, для которой $\phi(13) = 13.025 > 13$. Таким образом, для хорошо обусловленных систем этот вариант метода Гаусса является хорошо обусловленным.

Однако гарантия хорошей обусловленности достигается здесь ценой значительных затрат на выбор главных элементов. Для этого дополнительно к $(2/3)m^3$ арифметическим действиям требуется произвести примерно $m^3/3$ операций сравнения, что может ощутимо замедлить процесс решения задачи на компьютере². Поэтому в большинстве случаев на практике предпочтение отдается все же схеме частичного выбора. Как уже отмечено выше, ситуации, когда при использовании этого варианта метода Гаусса происходит существенный рост элементов,

¹ N. Gould. On growth in Gaussian elimination with complete pivoting // SIAM J. Matrix Anal. Appl. 1991. Vol. 12 (2). P. 354—361.

² Заметим, что на современных высокопроизводительных компьютерах операции с плавающей точкой выполняются почти так же быстро, как и операции сравнения.

встречаются чрезвычайно редко. Более того, эти ситуации могут быть легко выявлены с помощью заложенных в современных программах эффективных методов слежения за ростом элементов матриц.

4. Случай, когда выбор главных элементов не нужен. Известно что, для некоторых классов матриц при использовании схемы единственного деления главные элементы гарантированно располагаются на главной диагонали и потому применять частичный выбор нет необходимости. Так, например, обстоит дело для систем с симметричными положительно определенными матрицами, а также с матрицами, обладающими свойством *строчного диагонального преобладания*:

$$\sum_{\substack{j=1 \\ j \neq i}}^m |a_{ij}| < |a_{ii}|, \quad i = 1, 2, \dots, m. \quad (5.45)$$

Матрицы, удовлетворяющие условию (5.45), таковы, что в каждой из строк модуль элемента a_{ii} , расположенного на главной диагонали, больше суммы модулей всех остальных элементов строки.

5. Масштабирование. Перед началом решения целесообразно масштабировать систему так, чтобы ее коэффициенты были величинами порядка единицы.

Существуют два естественных способа масштабирования¹ системы $Ax = b$. Первый заключается в умножении каждого из уравнений на некоторый масштабирующий множитель μ_i . Второй состоит в умножении на масштабирующий множитель α_j каждого j -го столбца матрицы, что соответствует замене переменных $x'_j = \alpha_j^{-1} x_j$ (фактически — это замена единиц измерения). В реальных ситуациях чаще всего масштабирование может быть выполнено без существенных трудностей. Однако подчеркнем, что в общем случае удовлетворительного способа масштабирования пока не найдено.

На практике масштабирование обычно производят с помощью деления каждого уравнения на его наибольший по модулю коэффициент либо на евклидову норму соответствующей строки матрицы A . Это вполне удовлетворительный способ для большинства реально встречающихся задач.

¹ Иногда эту операцию называют уравновешиванием.

§ 5.6. Метод Гаусса и решение систем уравнений с несколькими правыми частями, обращение матриц, вычисление определителей

Рассмотрим применение метода Гаусса к решению следующих задач линейной алгебры:

- 1) вычисление решений системы уравнений с несколькими правыми частями;
- 2) вычисление обратной матрицы;
- 3) вычисление определителя.

1. Вычисление решений системы уравнений с несколькими правым частями. Довольно часто на практике встречается ситуация, когда нужно решить несколько систем уравнений

$$Ax = d_{(1)}, Ax = d_{(2)}, \dots, Ax = d_{(p)} \quad (5.46)$$

с одной матрицей A и различными правыми частями $d_{(1)}, d_{(2)}, \dots, d_{(p)}$.

Конечно, применив метод Гаусса к каждой из систем (5.46) независимо от других, можно найти соответствующие решения $x_{(1)}, x_{(2)}, \dots, x_{(p)}$, затратив примерно $(2/3)pm^3$ арифметических операций. Однако при одновременном решении систем (5.46) число операций можно существенно сократить. Как было отмечено в § 5.5, основные вычислительные затраты в методе Гаусса связаны с преобразованием матрицы к треугольному виду. Преобразование же правой части производится параллельно и требует примерно m^2 арифметических операций. Если параллельно с приведением матрицы A к треугольному виду преобразовать по однотипным формулам все p правых частей, то на прямой ход метода будет затрачено только примерно $(2/3)m^3 + pm^2$ операций. С учетом обратного хода общие вычислительные затраты составят примерно $(2/3)m^3 + 2pm^2$ арифметических операций.

Преобразование правых частей не обязательно производить параллельно. Каждую из правых частей $d_{(1)}, d_{(2)}, \dots, d_{(p)}$ можно обработать последовательно, если после приведения матрицы A к треугольному виду $A^{(m-1)}$ сохранить в памяти компьютера множители μ_{ij} и матрицу $A^{(m-1)}$.

2. Вычисление обратной матрицы. Прежде чем переходить к изложению метода вычисления обратной матрицы A^{-1} для квадратной невырожденной матрицы A , отметим, что в действительности проблема вычисления обратной матрицы возникает не так часто, как это можно предполагать.

К сожалению, нередко обращение матрицы A производится с единственной целью вычислить по известному вектору b вектор x вида $x = A^{-1}b$. Умножение матрицы A^{-1} на вектор требует примерно $2m^2$ арифметических операций. Однако вычисление A^{-1} обходится (как будет показано

ниже) примерно в $2m^3$ операций. Это означает, что на вычисление решения системы $\mathbf{Ax} = \mathbf{b}$ по формуле $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ будет затрачено примерно $2m^3$ операций. В данном случае \mathbf{x} можно найти в 3 раза быстрее методом Гаусса и вычисление \mathbf{A}^{-1} не нужно. Более того, можно ожидать, что вычисленное методом Гаусса решение окажется точнее, так как потребуется выполнение меньшего числа операций.

Может показаться особенно выгодным предварительное вычисление матрицы \mathbf{A}^{-1} , если далее потребуется найти большое число векторов по формулам

$$\mathbf{x}_{(1)} = \mathbf{A}^{-1}\mathbf{d}_{(1)}, \quad \mathbf{x}_{(2)} = \mathbf{A}^{-1}\mathbf{d}_{(2)}, \dots, \quad \mathbf{x}_{(p)} = \mathbf{A}^{-1}\mathbf{d}_{(p)}. \quad (5.47)$$

Однако суммарные затраты при таком подходе составят примерно $2m^3 + 2pm^2$ операций, в то время как при одновременном решении методом Гаусса эквивалентных систем (5.46) получаем значения $\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \dots, \mathbf{x}_{(p)}$ примерно за $(2/3)m^3 + 2pm^2$ операций. Следовательно, и в этом случае вычисление \mathbf{A}^{-1} нецелесообразно.

Иногда в пользу необходимости вычисления \mathbf{A}^{-1} приводится следующий довод. Если известно, что в течение длительного времени потребуется неоднократно решать системы уравнений вида (5.46) с фиксированной матрицей \mathbf{A} и различными правыми частями $\mathbf{d}_{(k)}$, то имеет смысл предварительно вычислить \mathbf{A}^{-1} . Записав \mathbf{A}^{-1} в память компьютера можно затем по мере необходимости быстро вычислить $\mathbf{x}_{(k)}$ по формуле $\mathbf{x}_{(k)} = \mathbf{A}^{-1}\mathbf{d}_{(k)}$. Однако использование LU -разложения матрицы \mathbf{A} (см. § 5.7) позволяет вычислять $\mathbf{x}_{(k)}$ столь же быстро, а предварительная работа на этапе разложения дает трехкратную экономию. Таким образом, и этот довод в пользу необходимости вычисления обратной матрицы неубедителен.

Довольно часто при решении различных задач средствами линейной алгебры возникают выражения типа

$$\mathbf{v} = \mathbf{B}^{-1}\mathbf{C}\mathbf{A}^{-1}\mathbf{W}\mathbf{D}^{-1}\mathbf{w}. \quad (5.48)$$

Если у исследователя нет достаточного опыта решения задач линейной алгебры на компьютере, то он может принять решение о необходимости вычислять матрицы $\mathbf{B}^{-1}, \mathbf{A}^{-1}, \mathbf{D}^{-1}$, чтобы действовать далее по формуле (5.48). Однако и в этом случае, можно поступить иначе и найти \mathbf{v} с меньшими затратами. Решив систему $\mathbf{Dx} = \mathbf{w}$, найдем $\mathbf{x} = \mathbf{D}^{-1}\mathbf{w}$. Затем вычислим $\mathbf{y} = \mathbf{Wx}$ и, решив систему $\mathbf{Az} = \mathbf{y}$, найдем $\mathbf{z} = \mathbf{A}^{-1}\mathbf{y}$. Наконец, вычислим $\mathbf{u} = \mathbf{Cz}$ и, решив систему $\mathbf{Bv} = \mathbf{u}$, найдем $\mathbf{v} = \mathbf{B}^{-1}\mathbf{u}$.

Сказанное выше вовсе не означает, что нет ситуаций, когда вычисление матрицы \mathbf{A}^{-1} необходимо и оправдано. В ряде технических приложений и статистических задач непосредственный интерес представляет анализ свойств именно обратной матрицы. Тем не менее, как мы

видим, в матричных вычислениях можно и следует обходиться без вычисления обратных матриц. Авторы настоятельно рекомендуют не вычислять обратные матрицы, если только в дальнейшем не предполагается анализ элементов этих матриц.

Покажем, как вычисление обратной матрицы можно свести к рассмотренной выше задаче решения системы уравнений с несколькими правыми частями. Обозначим матрицу A^{-1} через V , ее столбцы — через v_1, v_2, \dots, v_m и столбцы единичной матрицы E — через e_1, e_2, \dots, e_m .

Согласно определению обратной матрицы верно равенство $AV = E$, эквивалентное совокупности равенств

$$Av_1 = e_1, \quad Av_2 = e_2, \dots, \quad Av_m = e_m. \quad (5.49)$$

Таким образом, столбцы матрицы $V = A^{-1}$ (а следовательно, и саму матрицу), можно найти, решив m систем уравнений с общей матрицей A . Согласно изложенному выше, для этого потребовалось бы примерно $(8/3)m^3$ арифметических операций, однако учет специального вида правых частей системы (5.49) позволяет вычислять матрицу A^{-1} примерно за $2m^3$ операций.

3. Вычисление определителя. Воспользуемся алгоритмом метода Гаусса с выбором главного элемента по столбцу и заметим, что искомый определитель и определитель полученной треугольной матрицы $A^{(m-1)}$ связаны равенством

$$\det A = (-1)^s \det A^{(m-1)},$$

где s — число потребовавшихся перестановок строк. Остается воспользоваться формулой (5.17), и тогда получим

$$\det A = (-1)^s a_{11}^{(0)} a_{22}^{(1)} \dots a_{mm}^{(m-1)}, \quad (5.50)$$

где $a_{11}^{(0)} = a_{11}$. Отметим, что вычисление по формуле (5.50) требует особой аккуратности, в особенности если число m велико. Как мы убедились в примере 3.29, при вычислении произведений следует специальным образом упорядочивать сомножители. Неудачный порядок их расположения может привести к получению результата, равного $\pm\infty$, или к существенной потере точности.

Этого можно избежать, если для вычисления $\det A$ воспользоваться формулой

$$\ln |\det A| = \sum_{i=1}^m \ln |a_{ii}^{(i-1)}|.$$

Однако следует иметь в виду, что ее использование может привести к некоторой потере точности.

Замечание. Действительная необходимость в вычислении определителей возникает довольно редко. Во всяком случае основанные на их использовании алгоритмы оказываются весьма неэффективными (как в примере 3.34, где обсуждалось правило Крамера), и поэтому вычисление определителей давно уже не является элементом современных алгоритмов линейной алгебры. Кроме того, из результатов в § 5.4 следует, что использование величины $\det A$ для определения степени близости системы уравнений к вырожденной дает весьма ненадежный и сомнительный критерий.

Пример 5.11. Используя метод Гаусса с выбором главного элемента по столбцу, вычислим определитель матрицы

$$A = \begin{bmatrix} 2 & -9 & 5 \\ 1.2 & -5.3999 & 6 \\ 1 & -1 & -7.5 \end{bmatrix}$$

на 6-разрядном десятичном компьютере.

Повторив преобразования из примера 5.9, получим матрицу

$$A^{(2)} = \begin{bmatrix} 2 & -9 & 5 \\ 0 & 3.5 & -10 \\ 0 & 0 & 3.00029 \end{bmatrix}.$$

Так как была сделана одна перестановка строк, то формула (5.50) дает $\det(A) \approx (-1) \cdot 2 \cdot 3.5 \cdot 3.00029 \approx -21.0020$. ▲

Можно с достаточной степенью уверенности предположить, что во многих технических науках, где традиционно используются определители, в ближайшее время неизбежен переход к использованию других, более содержательных характеристик линейных моделей. Такими естественными характеристиками могут служить, например, собственные числа и собственные векторы матриц (см. гл. 8).

§ 5.7. Метод Гаусса и разложение матрицы на множители. LU-разложение

Вернемся еще раз к методу Гаусса, чтобы рассмотреть его с более общих позиций. Излагаемый ниже подход оказался чрезвычайно плодотворным и привел не только к более глубокому пониманию метода, но и позволил создать высокоеффективные машинные алгоритмы его реализации, а также рассмотреть другие точные методы с единой точки зрения.

Рассмотрим сначала простейший вариант метода Гаусса для решения системы линейных алгебраических уравнений

$$\mathbf{A}\mathbf{x} = \mathbf{b}. \quad (5.51)$$

1. Схема единственного деления и LU-разложение. При выполнении вычислений 1-го шага исключения по схеме единственного деления система уравнений приводится к виду

$$\mathbf{A}^{(1)}\mathbf{x} = \mathbf{b}^{(1)}, \quad (5.52)$$

где

$$\mathbf{A}^{(1)} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1m} \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & \dots & a_{2m}^{(1)} \\ 0 & a_{32}^{(1)} & a_{33}^{(1)} & \dots & a_{3m}^{(1)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & a_{m2}^{(1)} & a_{m3}^{(1)} & \dots & a_{mm}^{(1)} \end{bmatrix}, \quad \mathbf{b}^{(1)} = \begin{bmatrix} b_1 \\ b_2^{(1)} \\ b_3^{(1)} \\ \vdots \\ b_m^{(1)} \end{bmatrix},$$

а коэффициенты $a_{ij}^{(1)}$, $b_i^{(1)}$ ($i, j = 2, 3, \dots, m$) вычисляются по формулам (5.29), (5.31).

Введем матрицу

$$\mathbf{M}_1 = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ -\mu_{21} & 1 & 0 & \dots & 0 \\ -\mu_{31} & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ -\mu_{m1} & 0 & 0 & \dots & 1 \end{bmatrix}.$$

Как нетрудно проверить, справедливы равенства

$$\mathbf{A}^{(1)} = \mathbf{M}_1 \mathbf{A}, \quad \mathbf{b}^{(1)} = \mathbf{M}_1 \mathbf{b},$$

т.е. преобразование системы (5.51) к виду (5.52) эквивалентно умножению левой и правой частей системы на матрицу \mathbf{M}_1 .

Аналогично можно показать, что вычисления 2-го шага исключения приводят систему (5.52) к виду

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

где

$$\mathbf{A}^{(2)} = \mathbf{M}_2 \mathbf{A}^{(1)}, \quad \mathbf{b}^{(2)} = \mathbf{M}_2 \mathbf{b}^{(1)};$$

$$\mathbf{A}^{(2)} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1m} \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & \dots & a_{2m}^{(1)} \\ 0 & 0 & a_{33}^{(2)} & \dots & a_{3m}^{(2)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & a_{m3}^{(2)} & \dots & a_{mm}^{(2)} \end{bmatrix}, \quad \mathbf{M}_2 = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & -\mu_{32} & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & -\mu_{m2} & 0 & \dots & 1 \end{bmatrix}, \quad \mathbf{b}^{(2)} = \begin{bmatrix} b_1 \\ b_2^{(1)} \\ b_3^{(2)} \\ \vdots \\ b_m^{(2)} \end{bmatrix}.$$

После $(m-1)$ -го шага, завершающего прямой ход, система оказывается приведенной к виду

$$\mathbf{A}^{(m-1)} \mathbf{x} = \mathbf{b}^{(m-1)} \quad (5.53)$$

с верхней треугольной матрицей $\mathbf{A}^{(m-1)}$. Здесь

$$\mathbf{A}^{(m-1)} = \mathbf{M}_{m-1} \mathbf{A}^{m-2}, \quad \mathbf{b}^{(m-1)} = \mathbf{M}_{m-1} \mathbf{b}^{(m-2)},$$

где

$$\mathbf{A}^{(m-1)} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1m} \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & \dots & a_{2m}^{(1)} \\ 0 & 0 & a_{33}^{(2)} & \dots & a_{3m}^{(2)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & a_{mm}^{(m-1)} \end{bmatrix}, \quad \mathbf{M}_{m-1} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & 0 & \dots & -\mu_{m,m-1} & 1 \end{bmatrix},$$

$$\mathbf{b}^{(m-1)} = \begin{bmatrix} b_1 \\ b_2^{(1)} \\ b_3^{(2)} \\ \vdots \\ b_m^{(m-1)} \end{bmatrix}.$$

Заметим, что матрица $\mathbf{A}^{(m-1)}$ получена из матрицы \mathbf{A} последовательным умножением на $\mathbf{M}_1, \mathbf{M}_2, \dots, \mathbf{M}_{m-1}$:

$$\mathbf{A}^{(m-1)} = \mathbf{M}_{m-1} \dots \mathbf{M}_2 \mathbf{M}_1 \mathbf{A}. \quad (5.54)$$

Аналогично,

$$\mathbf{b}^{(m-1)} = \mathbf{M}_{m-1} \dots \mathbf{M}_2 \mathbf{M}_1 \mathbf{b}. \quad (5.55)$$

Из равенства (5.54) вытекает следующее представление:

$$\mathbf{A} = \mathbf{M}_1^{-1} \mathbf{M}_2^{-1} \dots \mathbf{M}_{m-1}^{-1} \mathbf{A}^{(m-1)}. \quad (5.56)$$

Как легко проверить,

$$\mathbf{M}_1^{-1} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ \mu_{21} & 1 & 0 & \dots & 0 \\ \mu_{31} & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ \mu_{m1} & 0 & 0 & \dots & 1 \end{bmatrix}, \quad \mathbf{M}_2^{-1} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & \mu_{32} & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \mu_{m2} & 0 & \dots & 1 \end{bmatrix}, \dots,$$

$$\mathbf{M}_{m-1}^{-1} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & 0 & \dots & \mu_{m,m-1} & 1 \end{bmatrix}.$$

Для этого достаточно перемножить матрицы \mathbf{M}_k^{-1} и \mathbf{M}_k ($k = 1, \dots, m-1$), в результате чего получится единичная матрица.

Введем обозначения $\mathbf{U} = \mathbf{A}^{(m-1)}$, $\mathbf{L} = \mathbf{M}_1^{-1} \mathbf{M}_2^{-1} \dots \mathbf{M}_{m-1}^{-1}$. Вычисляя матрицу \mathbf{L} , убеждаемся в том, что она имеет следующий вид:

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ \mu_{21} & 1 & 0 & \dots & 0 \\ \mu_{31} & \mu_{32} & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ \mu_{m1} & \mu_{m2} & \mu_{m3} & \dots & 1 \end{bmatrix}. \quad (5.57)$$

Тогда равенство (5.56) в новых обозначениях примет вид

$$\mathbf{A} = \mathbf{LU}. \quad (5.58)$$

Это и есть¹ LU-разложение матрицы \mathbf{A} — представление матрицы \mathbf{A} в виде произведения нижней треугольной матрицы \mathbf{L} и верхней треугольной матрицы \mathbf{U} .

Таким образом, прямой ход метода Гаусса без перестановок можно рассматривать как процесс вычисления LU-разложения матрицы системы, на k -м шаге которого определяются элементы k -го столбца матрицы \mathbf{L} и k -й строки матрицы \mathbf{U} .

Возможность LU-разложения обосновывается следующей теоремой.

¹ Обозначения треугольных матриц буквами \mathbf{L} и \mathbf{U} вызваны тем, что эти буквы являются начальными в английских словах lower — «нижний» и upper — «верхний».

Теорема 5.3. Если все главные миноры матрицы A отличны от нуля, то существуют единственная нижняя треугольная матрица L вида (5.57) и верхняя треугольная матрица U такие, что $A = LU$. ■

Структура матриц L и U позволяет организовывать компактное размещение элементов этих матриц в памяти компьютера по мере их вычисления. На k -м шаге исключения в области памяти, где первоначально располагалась матрица A , размещается матрица

$$\left[\begin{array}{cccccc} a_{11} & a_{12} & \dots & a_{1k} & a_{1,k+1} & \dots & a_{1m} \\ \mu_{21} & a_{22}^{(1)} & \dots & a_{2k}^{(1)} & a_{2,k+1}^{(1)} & \dots & a_{2m}^{(1)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \mu_{k1} & \mu_{k2} & \dots & a_{kk}^{(k-1)} & a_{k,k+1}^{(k-1)} & \dots & a_{km}^{(k-1)} \\ \mu_{k+1,1} & \mu_{k+1,2} & \dots & \mu_{k+1,k} & a_{k+1,k+1}^{(k)} & \dots & a_{k+1,m}^{(k)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \mu_{m1} & \mu_{m2} & \dots & \mu_{mk} & a_{m,k+1}^{(k)} & \dots & a_{mm}^{(k)} \end{array} \right].$$

При этом вся необходимая для дальнейших вычислений информация сохраняется.

Пример 5.12. Проиллюстрируем LU -разложение на примере решения системы (5.35). На основании данных табл. 5.2 можно записать

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0.5 & 1 & 0 & 0 \\ 0.3 & -1.6 & 1 & 0 \\ 0 & -3 & 2.5 & 1 \end{bmatrix}, \quad U = \begin{bmatrix} 10 & 6 & 2 & 0 \\ 0 & -2 & -3 & 4 \\ 0 & 0 & -4.4 & 5.4 \\ 0 & 0 & 0 & 0.5 \end{bmatrix}.$$

Следовательно, LU -разложение матрицы системы имеет вид

$$\begin{bmatrix} 10 & 6 & 2 & 0 \\ 5 & 1 & -2 & 4 \\ 3 & 5 & 1 & -1 \\ 0 & 6 & -2 & .2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0.5 & 1 & 0 & 0 \\ 0.3 & -1.6 & 1 & 0 \\ 0 & -3 & 2.5 & 1 \end{bmatrix} \cdot \begin{bmatrix} 10 & 6 & 2 & 0 \\ 0 & -2 & -3 & 4 \\ 0 & 0 & -4.4 & 5.4 \\ 0 & 0 & 0 & 0.5 \end{bmatrix}. \blacktriangle$$

2. Использование LU -разложения. В современных программах, реализующих метод Гаусса на компьютере, вычисления разбивают на два основных этапа. Первый этап — это вычисление LU -разложения матрицы системы. Второй этап — обработка правых частей и вычисление решения.

Смысл выделения первого этапа состоит в том, что он может быть выполнен независимо, для его проведения не нужна информация о правой части системы. Это как бы этап предварительной подготовки к быстрому вычислению решения. Именно для получения LU-разложения производится основная масса вычислений (примерно $(2/3)m^3$ арифметических операций).

На втором этапе выполняют следующие действия:

1°. Преобразуют правую часть \mathbf{b} по формулам прямого хода; необходимые для вычисления коэффициенты μ_{ij} берут из матрицы L . В результате получают вектор $\mathbf{b}^{(m-1)}$, связанный с вектором \mathbf{b} формулой (5.55).

2°. С помощью обратной подстановки решают систему $U\mathbf{x} = \mathbf{b}^{(m-1)}$.

Для непосредственного вычисления решения \mathbf{x} на втором этапе требуется примерно $2m^2$ арифметических операций.

В случае, если необходимо решить p систем уравнений (5.46) с фиксированной матрицей A и различными правыми частями $\mathbf{d}_{(1)}, \mathbf{d}_{(2)}, \dots, \mathbf{d}_{(p)}$, первый этап проводят лишь один раз. Затем последовательно p раз проводят вычисления второго этапа для получения решений $\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \dots, \mathbf{x}_{(p)}$. Для этого, как и в § 5.6, требуется примерно $(2/3)m^2 + 2pm^2$ арифметических операций.

Пример 5.13. Решим систему

$$\begin{aligned} 10x_1 + 6x_2 + 2x_3 &= 8, \\ 5x_1 + x_2 - 2x_3 + 4x_4 &= 7, \\ 3x_1 + 5x_2 + x_3 - x_4 &= 2, \\ 6x_2 - 2x_3 + 2x_4 &= 2. \end{aligned}$$

Воспользуемся LU-разложением матрицы системы, указанным в примере 5.12. Сначала преобразуем вектор правой части $\mathbf{b} = (8, 7, 2, 2)^T$ по формулам прямого хода.

1-й шаг. $b_2^{(1)} = b_2 - \mu_{21}b_1 = 7 - 0.5 \cdot 8 = 3$, $b_3^{(1)} = b_3 - \mu_{31}b_1 = 2 - -0.3 \cdot 8 = -0.4$, $b_4^{(1)} = b_4 - \mu_{41}b_1 = 2 - 0 \cdot 8 = 2$. После 1-го шага получим $\mathbf{b}^{(1)} = (8, 3, -0.4, 2)^T$.

2-й шаг. $b_3^{(2)} = b_3^{(1)} - \mu_{32}b_2^{(1)} = -0.4 - (-1.6) \cdot 3 = 4.4$, $b_4^{(2)} = b_4^{(1)} - \mu_{42}b_2^{(1)} = 2 - (-3) \cdot 3 = 11$. После 2-го шага найдем $\mathbf{b}^{(2)} = (8, 3, 4.4, 11)^T$.

3-й шаг. $b_4^{(3)} = b_4^{(2)} - \mu_{43}b_3^{(2)} = 11 - 2.5 \cdot 4.4 = 0$. В результате прямого хода получен вектор $\mathbf{b}^{(3)} = (8, 3, 4.4, 0)^T$ и система приведена к виду

$$\begin{aligned} 10x_1 + 6x_2 + 2x_3 &= 8, \\ -2x_2 - 3x_3 + 4x_4 &= 3, \\ -4.4x_3 + 5.4x_4 &= 4.4, \\ 0.5x_4 &= 0. \end{aligned}$$

Обратный ход дает значения неизвестных $x_4 = 0, x_3 = -1, x_2 = 0, x_1 = 1$. ▲

3. Метод Гаусса с выбором главного элемента и разложение матрицы на множители. В отличие от схемы единственного деления схема частичного выбора предполагает на k -м шаге прямого хода перестановку уравнений системы с номерами i_k и k (при выборе в качестве главного элемента k -го шага элемента $a_{i_k k}^{(k-1)}$). Это преобразование эквивалентно умножению системы на матрицу P_k , которая получается из единичной матрицы перестановкой i_k -й и k -й строк (см. пример 5.14). Исключение неизвестного на k -м шаге по-прежнему эквивалентно умножению системы на матрицу M_k .

Таким образом, после 1-го шага система $Ax = b$ преобразуется к виду $A^{(1)}x = b^{(1)}$, где $A^{(1)} = M_1 P_1 A$, $b^{(1)} = M_1 P_1 b$. После 2-го шага система преобразуется к виду $A^{(2)}x = b^{(2)}$, где $A^{(2)} = M_2 P_2 A^{(1)}$, $b^{(2)} = M_2 P_2 b^{(1)}$.

После завершающего $(m-1)$ -го шага прямого хода система оказывается приведенной к виду $A^{(m-1)}x = b^{(m-1)}$, где $A^{(m-1)} = M_{m-1} P_{m-1} A^{(m-2)}$, $b^{(m-1)} = M_{m-1} P_{m-1} b^{(m-2)}$.

Как нетрудно видеть,

$$A^{(m-1)} = M_{m-1} P_{m-1} \dots M_2 P_2 M_1 P_1 A, \quad (5.59)$$

$$b^{(m-1)} = M_{m-1} P_{m-1} \dots M_2 P_2 M_1 P_1 b. \quad (5.60)$$

Равенство (5.59) равносильно следующему разложению матрицы A на множители:

$$A = P_1^{-1} M_1^{-1} P_2^{-1} M_2^{-1} \dots P_{m-1}^{-1} M_{m-1}^{-1} U, \quad (5.61)$$

где $U = A^{(m-1)}$ — верхняя треугольная матрица.

Разложение (5.61) не является LU -разложением матрицы A . Однако прямой ход по-прежнему равносителен LU -разложению, но уже не самой

матрицы A , а матрицы \tilde{A} , полученной из нее в результате соответствующей перестановки строк. Это разложение имеет вид

$$\tilde{A} = \tilde{L}U, \quad (5.62)$$

где $\tilde{A} = P_{m-1}P_{m-2}\dots P_2P_1A$, \tilde{L} — нижняя треугольная матрица, отличающаяся от матрицы (5.57) перестановкой множителей в столбцах.

Пример 5.14. Найдем разложение вида (5.62) для матрицы системы (5.39), используя результаты вычислений примера 5.9. Так как 1-й шаг прямого хода не потребовал перестановки, а на 2-м шаге были переставлены второе и третье уравнения, то

$$P_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad \tilde{A} = \begin{bmatrix} 2 & -9 & 5 \\ 1 & -1 & -7.5 \\ 1.2 & -5.3999 & 6 \end{bmatrix}.$$

Для матрицы \tilde{A} прямой ход уже проводится по схеме единственного деления. Отличие от вычислений примера 5.9 состоит в том, что на 2-м шаге множители μ_{21} и μ_{31} , а также второе и третье уравнения системы (5.40) меняются местами. В результате получим разложение вида (5.62), где

$$\tilde{L} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0.6 & 2.9 \cdot 10^{-5} & 1 \end{bmatrix}, \quad U = \begin{bmatrix} 2 & -9 & 5 \\ 0 & 3.5 & -10 \\ 0 & 0 & 3.00029 \end{bmatrix}. \blacksquare$$

После получения разложения вида (5.62) для решения системы $Ax = b$ выполняют следующие действия.

1°. Правую часть перестановкой элементов приводят к виду

$$\tilde{b} = P_{m-1}P_{m-2}\dots P_2P_1b.$$

2°. Преобразуют вектор \tilde{b} по формулам прямого хода; необходимые для вычислений множители $\tilde{\mu}_{ij}$ берут из матрицы \tilde{L} . В результате получают вектор $b^{(m-1)}$.

3°. Обратной подстановкой решают систему $Ux = b^{(m-1)}$.

Заметим, что матрица перестановки P_k полностью определяется заданием номера i_k уравнения, которое переставляется с k -м уравнением. Поэтому для хранения всей информации о перестановках достаточно целочисленного массива длины $m - 1$.

Пример 5.15. Решим систему (5.39), используя полученное в примере 5.14 разложение матрицы \tilde{A} .

Здесь вектор $b = (-4, 0.6001, -8.5)^T$ преобразуется в вектор $\tilde{b} = (-4, -8.5, 0.6001)^T$ перестановкой второго и третьего элементов.

Прямой ход.

1-й шаг. $b_2^{(1)} = \tilde{b}_2 - \tilde{\mu}_{21}\tilde{b}_1 = -8.5 - 0.5(-4) = -6.5$, $b_3^{(1)} = \tilde{b}_3 - \tilde{\mu}_{31}\tilde{b}_1 = 0.6001 - 0.6 \cdot (-4) = 3.0001$. В результате 1-го шага имеем $b^{(1)} = (-4, -6.5, 3.0001)^T$.

2-й шаг. $b_3^{(2)} = b_3^{(1)} - \tilde{\mu}_{32}b_2^{(1)} = 3.0001 - 2.9 \cdot 10^{-5} \cdot (-6.5) \approx 3.00029$.

В результате прямого хода правая часть оказалась приведенной к виду $b^{(2)} = (-4, -6.5, 3.00029)^T$.

Обратный ход проводится точно так же, как в примере 5.9, и дает значения $x_3 = 1$, $x_2 = 1$, $x_1 = 0$. ▲

§ 5.8. Метод Холецкого (метод квадратных корней)

1. Описание метода. Пусть требуется решить систему линейных алгебраических уравнений

$$Ax = b \quad (5.63)$$

с симметричной положительно определенной матрицей A . Системы уравнений такого типа часто встречаются в приложениях (например, в задачах оптимизации, при решении уравнений математической физики и др.). Для их решения весьма часто применяется *метод Холецкого*¹ (другое название — *метод квадратных корней*)².

В основе метода лежит алгоритм построения специального *LU*-разложения матрицы A , в результате чего она приводится к виду

$$A = LL^T. \quad (5.64)$$

В *разложении Холецкого* (5.64) нижняя треугольная матрица

$$\mathbf{L} = \begin{bmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ l_{m1} & l_{m2} & \dots & l_{mm} \end{bmatrix} \quad (5.65)$$

¹ Андре-Луи Холецкий (1875—1918) — французский математик и геодезист, лейтенант артиллерии. Погиб во время Первой мировой войны.

² Идея метода была опубликована Гауссом в 1806 г.

уже не обязательно должна иметь на главной диагонали единицы, как это было в методе Гаусса, а требуется только, чтобы диагональные элементы l_{ii} были положительными.

Если разложение (5.64) получено, то решение системы (5.63) сводится к последовательному решению двух систем с треугольными матрицами:

$$Ly = b, \quad L^T x = y. \quad (5.66)$$

Для решения систем (5.66) требуется выполнение примерно $2m^2$ арифметических операций.

Найдем элементы матрицы L . Для этого вычислим элементы матрицы LL^T и приравняем их соответствующим элементам матрицы A . В результате получим систему уравнений

$$\begin{aligned}
 l_{11}^2 &= a_{11}; \\
 l_{i1}l_{11} &= a_{i1}, \quad i = 2, 3, \dots, m; \\
 l_{21}^2 + l_{22}^2 &= a_{22}; \\
 l_{i1}l_{21} + l_{i2}l_{22} &= a_{i2}, \quad i = 3, 4, \dots, m; \\
 \dots & \\
 l_{k1}^2 + l_{k2}^2 + \dots + l_{kk}^2 &= a_{kk}; \\
 l_{i1}l_{k1} + l_{i2}l_{k2} + \dots + l_{ik}l_{kk} &= a_{ik}, \quad i = k+1, \dots, m; \\
 \dots & \\
 l_{m1}^2 + l_{m2}^2 + \dots + l_{mm}^2 &= a_{mm}.
 \end{aligned} \tag{5.67}$$

Решая систему (5.67), последовательно находим

$$\begin{aligned}
 l_{11} &= \sqrt{a_{11}}; \\
 l_{i1} &= a_{i1}/l_{11}, \quad i = 2, 3, \dots, m; \\
 l_{22} &= \sqrt{a_{22} - l_{21}^2}; \\
 l_{i2} &= (a_{i2} - l_{i1}l_{21})/l_{22}, \quad i = 3, 4, \dots, m; \\
 l_{kk} &= \sqrt{a_{kk} - l_{k1}^2 - l_{k2}^2 - \dots - l_{k,k-1}^2}; \\
 l_{ik} &= (a_{ik} - l_{i1}l_{k1} - l_{i2}l_{k2} - \dots - l_{i,k-1}l_{k,k-1})/l_{kk}, \quad i = k+1, \dots, m; \\
 l_{mm} &= \sqrt{a_{mm} - l_{m1}^2 - l_{m2}^2 - \dots - l_{m,m-1}^2}.
 \end{aligned} \tag{5.68}$$

Заметим, что для вычисления диагональных элементов используется операция извлечения квадратного корня. Поэтому метод Холецкого называют еще и методом квадратных корней. Доказано, что положительность соответствующих подкоренных выражений является следствием положительной определенности матрицы A .

Замечание. Как следует из формул (5.68), матрица L , входящая в разложение Холецкого (5.64), определяется по матрице A однозначно.

2. Достоинства метода. Метод Холецкого обладает рядом ценных качеств, которые позволяют предпочесть его методу Гаусса, если требуется решить систему линейных алгебраических уравнений с симметричной и положительно определенной матрицей.

Как нетрудно подсчитать, число операций, выполняемых в ходе вычисления разложения (5.64) по формулам (5.68), равно примерно $m^3/3$. Учитывая, что для решения систем (5.66) требуется примерно $2m^2$ арифметических операций, убеждаемся, что при больших m метод Холецкого требует вдвое меньше вычислительных затрат по сравнению с методом Гаусса.

Учет симметричности матрицы A позволяет экономно использовать память компьютера при записи исходных данных задачи и результатов вычислений. Действительно, для задания матрицы A достаточно ввести в память компьютера только элементы a_{ij} ($i \geq j$), расположенные на главной диагонали и под ней. В формулах (5.68) каждый такой элемент a_{ij} используется лишь однажды для получения l_{ij} и далее в вычислениях не участвует. Поэтому в процессе вычислений найденные элементы l_{ij} могут последовательно замещать элементы a_{ij} .

В результате нижняя треугольная матрица L может быть расположена в той области памяти, где первоначально хранилась нижняя треугольная часть матрицы A . Применение для решения системы (5.63) метода Гаусса потребовало бы использования примерно вдвое большего объема памяти.

Безусловным достоинством метода Холецкого является также его гарантированная устойчивость.

Пример 5.16. Используя метод Холецкого, найдем решение системы уравнений с симметричной положительно определенной матрицей:

$$\begin{aligned} 6.25x_1 - & x_2 + 0.5x_3 = 7.5, \\ -x_1 + & 5x_2 + 2.12x_3 = -8.68, \\ 0.5x_1 + 2.12x_2 + & 3.6x_3 = -0.24. \end{aligned}$$

По формулам (5.68) последовательно находим:

$$l_{11} = \sqrt{a_{11}} = \sqrt{6.25} = 2.5, \quad l_{21} = a_{21}/l_{11} = -1/2.5 = -0.4;$$

$$l_{31} = a_{31}/l_{11} = 0.5/2.5 = 0.2, \quad l_{22} = \sqrt{a_{22} - l_{21}^2} = \sqrt{5 - 0.16} = 2.2;$$

$$l_{32} = (a_{32} - l_{31}l_{21})/l_{22} = (2.12 - 0.2(-0.4))/2.2 = 1;$$

$$l_{33} = \sqrt{a_{33} - l_{31}^2 - l_{32}^2} = \sqrt{3.6 - 0.2^2 - 1^2} = 1.6.$$

Следовательно, матрица L такова:

$$L = \begin{bmatrix} 2.5 & 0 & 0 \\ -0.4 & 2.2 & 0 \\ 0.2 & 1 & 1.6 \end{bmatrix}.$$

Система $Ly = b$ имеет вид

$$2.5y_1 = 7.5,$$

$$-0.4y_1 + 2.2y_2 = -8.68,$$

$$0.2y_1 + y_2 + 1.6y_3 = -0.24.$$

Решая ее, получаем $y_1 = 3$, $y_2 = -3.4$, $y_3 = 1.6$.

Далее из системы $L^T x = y$, которая имеет вид

$$2.5x_1 - 0.4x_2 + 0.2x_3 = 3,$$

$$2.2x_2 + x_3 = -3.4,$$

$$1.6x_3 = 1.6,$$

находим решение $x_1 = 0.8$, $x_2 = -2$, $x_3 = 1$. \blacktriangle

§ 5.9. Метод прогонки

Рассмотрим *метод прогонки*¹ — простой и эффективный алгоритм решения систем линейных алгебраических уравнений с трехдиагональными матрицами

$$\begin{aligned} b_1x_1 + c_1x_2 &= d_1, \\ a_2x_1 + b_2x_2 + c_2x_3 &= d_2, \\ \dots &\dots \\ a_ix_{i-1} + b_ix_i + c_ix_{i+1} &= d_i, \\ \dots &\dots \\ a_{m-1}x_{m-2} + b_{m-1}x_{m-1} + c_{m-1}x_m &= d_{m-1}, \\ a_mx_{m-1} + b_mx_m &= d_m. \end{aligned} \tag{5.69}$$

¹ Метод прогонки был предложен в начале 50-х годов независимо несколькими авторами, в том числе российскими учеными И.М. Гельфандом, О.В. Локуциевским, В.С. Владимировым, А.С. Кронродом.

Системы такого вида часто возникают при решении различных задач математической физики, а также при решении других вычислительных задач (например, приближения функций сплайнами).

1. Вывод расчетных формул. Преобразуем первое уравнение системы (5.69) к виду

$$x_1 = \alpha_1 x_2 + \beta_1, \quad (5.70)$$

где $\alpha_1 = -c_1/b_1$, $\beta_1 = d_1/b_1$.

Подставим полученное для x_1 выражение во второе уравнение системы:

$$a_2(\alpha_1 x_2 + \beta_1) + b_2 x_2 + c_2 x_3 = d_2.$$

Преобразуем это уравнение к виду

$$x_2 = \alpha_2 x_3 + \beta_2, \quad (5.71)$$

где $\alpha_2 = -c_2/(b_2 + a_2\alpha_1)$, $\beta_2 = (d_2 - a_2\beta_1)/(b_2 + a_2\alpha_1)$.

Выражение (5.71) подставляем в третье уравнение системы и т.д.

На i -м шаге этого процесса ($1 < i < m$) i -е уравнение системы преобразуется к виду

$$x_i = \alpha_i x_{i+1} + \beta_i, \quad (5.72)$$

где $\alpha_i = -c_i/(b_i + a_i\alpha_{i-1})$, $\beta_i = (d_i - a_i\beta_{i-1})/(b_i + a_i\alpha_{i-1})$.

На m -м шаге подстановка в последнее уравнение выражения $x_{m-1} = \alpha_{m-1} x_m + \beta_{m-1}$ дает:

$$a_m(\alpha_{m-1} x_m + \beta_{m-1}) + b_m x_m = d_m.$$

Откуда можно определить значение x_m :

$$x_m = \beta_m = (d_m - a_m \beta_{m-1})/(b_m + a_m \alpha_{m-1}).$$

Значения остальных неизвестных x_i для $i = m-1, m-2, \dots, 1$ теперь легко вычисляются по формуле (5.72).

2. Алгоритм прогонки. Сделанные преобразования позволяют организовать вычисления метода прогонки в два этапа.

Прямой ход метода прогонки (*прямая прогонка*) состоит в вычислении прогоночных коэффициентов α_i ($1 \leq i < m$) и β_i ($1 \leq i \leq m$). При $i = 1$ коэффициенты вычисляют по формулам

$$\alpha_1 = -c_1/\gamma_1, \quad \beta_1 = d_1/\gamma_1, \quad \gamma_1 = b_1, \quad (5.73)$$

а при $i = 2, 3, \dots, m-1$ — по рекуррентным формулам

$$\alpha_i = -c_i/\gamma_i, \quad \beta_i = (d_i - a_i \beta_{i-1})/\gamma_i, \quad \gamma_i = b_i + a_i \alpha_{i-1} \quad (5.74)$$

При $i = m$ прямая прогонка завершается вычислением

$$\beta_m = (d_m - a_m \beta_{m-1}) / \gamma_m, \quad \gamma_m = b_m + a_m \alpha_{m-1}. \quad (5.75)$$

Обратный ход метода прогонки (*обратная прогонка*), дает значения неизвестных. Сначала полагают $x_m = \beta_m$. Затем значения остальных неизвестных вычисляют по формуле

$$x_i = \alpha_i x_{i+1} + \beta_i, \quad i = m-1, m-2, \dots, 1. \quad (5.76)$$

Вычисления ведут в порядке убывания значений i от $m-1$ до 1.

Пример 5.17. Используя метод прогонки, решим систему

$$5x_1 - x_2 = 2.0,$$

$$2x_1 + 4.6x_2 - x_3 = 3.3,$$

$$2x_2 + 3.6x_3 - 0.8x_4 = 2.6,$$

$$3x_3 + 4.4x_4 = 7.2.$$

Прямой ход. Согласно формулам (5.73)–(5.75) получаем:

$$\gamma_1 = b_1 = 5, \quad \alpha_1 = -c_1/\gamma_1 = 1/5 = 0.2, \quad \beta_1 = d_1/\gamma_1 = 2.0/5 = 0.4;$$

$$\gamma_2 = b_2 + a_2 \alpha_1 = 4.6 + 2 \cdot 0.2 = 5, \quad \alpha_2 = -c_2/\gamma_2 = 1/5 = 0.2;$$

$$\beta_2 = (d_2 - a_2 \beta_1)/\gamma_2 = (3.3 - 2 \cdot 0.4)/5 = 0.5;$$

$$\gamma_3 = b_3 + a_3 \alpha_2 = 3.6 + 2 \cdot 0.2 = 4, \quad \alpha_3 = -c_3/\gamma_3 = 0.8/4 = 0.2;$$

$$\beta_3 = (d_3 - a_3 \beta_2)/\gamma_3 = (2.6 - 2 \cdot 0.5)/4 = 0.4;$$

$$\gamma_4 = b_4 + a_4 \alpha_3 = 4.4 + 3 \cdot 0.2 = 5;$$

$$\beta_4 = (d_4 - a_4 \beta_3)/\gamma_4 = (7.2 - 3 \cdot 0.4)/5 = 1.2.$$

Обратный ход. Полагаем $x_4 = \beta_4 = 1.2$. Далее находим:

$$x_3 = \alpha_3 x_4 + \beta_3 = 0.2 \cdot 1.2 + 0.4 = 0.64;$$

$$x_2 = \alpha_2 x_3 + \beta_2 = 0.2 \cdot 0.64 + 0.5 = 0.628;$$

$$x_1 = \alpha_1 x_2 + \beta_1 = 0.2 \cdot 0.628 + 0.4 = 0.5256.$$

Итак, получаем решение:

$$x_1 = 0.5256, \quad x_2 = 0.628, \quad x_3 = 0.64, \quad x_4 = 1.2. \quad \blacktriangle$$

3. Свойства метода прогонки. Непосредственный подсчет показывает, что для реализации вычислений по формулам (5.73)–(5.76) требуется примерно $8m$ арифметических операций, тогда как в методе Гаусса это число составляет примерно $(2/3)m^3$. Важно и то, что трехдиагональная структура матрицы системы позволяет использовать для ее хранения лишь $3m - 2$ машинных слова.

Таким образом, при одной и той же производительности и оперативной памяти компьютера метод прогонки позволяет решать системы

гораздо большей размерности, чем стандартный метод Гаусса для систем уравнений с заполненной матрицей.

Приведем простые достаточные условия на коэффициенты системы (5.69), при выполнении которых вычисления по формулам прямой прогонки могут быть доведены до конца (ни один из знаменателей γ_i не обратится в нуль). В частности, это гарантирует существование решения системы (5.69) и его единственность. При выполнении тех же условий коэффициенты a_i , при всех i удовлетворяют неравенству $|a_i| \leq 1$, а следовательно, обратная прогонка по формуле (5.76) устойчива по входным данным (см. пример 3.27). Положим $a_1 = 0$, $c_m = 0$.

Теорема 5.4. Пусть коэффициенты системы (5.69) удовлетворяют следующим условиям диагонального преобладания:

$$|b_k| \geq |a_k| + |c_k|, \quad |b_k| > |a_k| \quad (1 \leq k \leq m). \quad (5.77)$$

Тогда $\gamma_i \neq 0$ и $|\alpha_i| \leq 1$ для всех $i = 1, 2, \dots, m$.

Доказательство. Воспользуемся принципом математической индукции. Из условий теоремы имеем $\gamma_1 = b_1 \neq 0$ и $|\alpha_1| = |c_1| / |\gamma_1| \leq 1$.

Пусть теперь $\gamma_{k-1} \neq 0$ и $|\alpha_{k-1}| \leq 1$ для некоторого $k > 1$. Тогда

$$|\gamma_k| = |b_k + a_k \alpha_{k-1}| \geq |b_k| - |a_k| |\alpha_{k-1}| \geq |b_k| - |a_k|.$$

Из полученной оценки в силу условий (5.77) вытекает справедливость неравенств $|\gamma_k| > 0$ и $|\gamma_k| \geq |c_k|$. Следовательно, $\gamma_k \neq 0$ и $|\alpha_k| = |c_k| / |\gamma_k| \leq 1$. ■

4. Метод прогонки и разложение матрицы на множители. Описанный вариант метода прогонки можно рассматривать как одну из схем метода Гаусса (без выбора главного элемента), в результате прямого хода которого исходная трехдиагональная матрица

$$\mathbf{A} = \begin{bmatrix} b_1 & c_1 & 0 & 0 & \dots & 0 & 0 & 0 \\ a_2 & b_2 & c_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & a_3 & b_3 & c_3 & \dots & 0 & 0 & 0 \\ \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & a_{m-1} & b_{m-1} & c_{m-1} \\ 0 & 0 & 0 & 0 & \dots & 0 & a_m & b_m \end{bmatrix}$$

представляется в виде произведения двух двухдиагональных матриц:

$$\mathbf{A} = \mathbf{L}\mathbf{U}. \quad (5.78)$$

Здесь

$$\mathbf{L} = \begin{bmatrix} \gamma_1 & 0 & 0 & \dots & 0 & 0 \\ a_2 & \gamma_2 & 0 & \dots & 0 & 0 \\ 0 & a_3 & \gamma_3 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & a_m & \gamma_m \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} 1 & -\alpha_1 & 0 & \dots & 0 & 0 \\ 0 & 1 & -\alpha_2 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 1 & -\alpha_{m-1} \\ 0 & 0 & 0 & \dots & 0 & 1 \end{bmatrix}.$$

Так как для определения \mathbf{L} и \mathbf{U} нет необходимости вычислять коэффициенты β_i , то общее число операций на получение разложения (5.78) составляет примерно $3m$.

Подобно тому как это описано в § 5.7, разложение (5.78) можно использовать для решения систем с многими правыми частями. Если нужно решить p систем с матрицей \mathbf{A} , то общее число операций составит примерно $3mp + 5mp$.

К сожалению, при обращении матрицы \mathbf{A} теряется ее трехдиагональная структура. Обратная матрица является заполненной, однако для ее вычисления с помощью разложения (5.78) требуется примерно лишь $2.5m^2$ арифметических операций.

Так как $\det \mathbf{A} = \det \mathbf{L} \cdot \det \mathbf{U}$, а $\det \mathbf{U} = 1$, то определитель трехдиагональной матрицы, после того как получено разложение (5.78), вычисляется по элементарной формуле

$$\det \mathbf{A} = \gamma_1 \gamma_2 \dots \gamma_m.$$

5. Некоторые варианты метода прогонки. Наряду с изложенным выше «стандартным» вариантом метода прогонки (правой прогонкой) существует большое число других вариантов этого метода. Это методы левой прогонки, встречных прогонок, немонотонной прогонки, потоковый вариант метода прогонки. В ряде случаев эти модификации могут существенно улучшить обусловленность прогонки.

Для систем уравнений, обладающих близкой к (5.69) структуре, разработаны методы циклической прогонки, матричной прогонки и др.

С указанными вариантами метода прогонки можно подробно ознакомиться в [50, 89].

§ 5.10. *QR*-разложение матрицы. Методы вращений и отражений

Метод Гаусса не является единственным методом исключения, используемым для решения систем линейных уравнений и приведения матриц к треугольному виду. Рассмотрим два метода исключения, обладающих в отличие от метода Гаусса гарантированной хорошей обусловленностью — метод вращений и метод отражений. Оба этих метода позволяют получить представление исходной матрицы A в виде произведения ортогональной¹ матрицы Q на верхнюю треугольную матрицу R :

$$A = QR. \quad (5.79)$$

Представление (5.79) — это *QR*-разложение матрицы на множители

1. Метод вращений. Опишем прямой ход метода. На 1-м шаге неизвестное x_1 исключают из всех уравнений, кроме первого. Для исключения x_1 из второго уравнения вычисляют числа

$$c_{12} = \frac{a_{11}}{\sqrt{a_{11}^2 + a_{21}^2}}, \quad s_{12} = \frac{a_{21}}{\sqrt{a_{11}^2 + a_{21}^2}}, \quad (5.80)$$

обладающие следующими свойствами:

$$c_{12}^2 + s_{12}^2 = 1, \quad -s_{12}a_{11} + c_{12}a_{21} = 0. \quad (5.81)$$

Затем первое уравнение системы заменяют линейной комбинацией первого и второго уравнений с коэффициентами c_{12} и s_{12} , а второе уравнение — аналогичной линейной комбинацией с коэффициентами $-s_{12}$ и c_{12} . В результате получают систему

$$\begin{aligned} a_{11}^{(1)}x_1 + a_{12}^{(1)}x_2 + a_{13}^{(1)}x_3 + \dots + a_{1m}^{(1)}x_m &= b_1^{(1)}, \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2m}^{(1)}x_m &= b_2^{(1)}, \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 + \dots + a_{3m}x_m &= b_3, \\ \dots \\ a_{m1}x_1 + a_{m2}x_2 + a_{m3}x_3 + \dots + a_{mm}x_m &= b_m, \end{aligned} \quad (5.82)$$

¹ Напомним, что вещественная матрица Q называется ортогональной, если для нее выполнено условие $Q^T = Q^{-1}$, что эквивалентно равенствам $Q^T Q = Q Q^T = E$. Важное свойство ортогонального преобразования векторов (т.е. преобразования векторов с помощью их умножения на ортогональную матрицу Q), состоит в том, что это преобразование не меняет евклидову норму векторов $\|Qx\|_2 = \|x\|_2$ для всех $x \in \mathbb{R}^n$.

в которой

$$\begin{aligned} a_{1j}^{(1)} &= c_{12}a_{1j} + s_{12}a_{2j}, \quad a_{2j}^{(1)} = -s_{12}a_{1j} + c_{12}a_{2j}, \quad (1 \leq j \leq m); \\ b_1^{(1)} &= c_{12}b_1 + s_{12}b_2, \quad b_2^{(1)} = -s_{12}b_1 + c_{12}b_2 \end{aligned} \quad (5.83)$$

Заметим, что $a_{21}^{(1)} = -s_{12}a_{11} + c_{12}a_{21} = 0$ в силу специального выбора чисел c_{12} и s_{12} (см. (5.81)).

Естественно, что если $a_{21} = 0$, исходная система уже имеет вид (5.82) и в исключении неизвестного x_1 из второго уравнения нет необходимости. В этом случае полагают $c_{12} = 1$ и $s_{12} = 0$.

Как нетрудно видеть, преобразование исходной системы (5.1) к виду (5.82) эквивалентно умножению слева матрицы A и правой части b на матрицу T_{12} , имеющую вид

$$T_{12} = \begin{bmatrix} c_{12} & s_{12} & 0 & 0 & \dots & 0 \\ -s_{12} & c_{12} & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 1 \end{bmatrix}.$$

Для исключения неизвестного x_1 из третьего уравнения вычисляют числа

$$c_{13} = \frac{a_{11}^{(1)}}{\sqrt{(a_{11}^{(1)})^2 + a_{31}^2}}, \quad s_{13} = \frac{a_{31}}{\sqrt{(a_{11}^{(1)})^2 + a_{31}^2}} \quad (5.84)$$

такие, что $c_{13}^2 + s_{13}^2 = 1$, $-s_{13}a_{11}^{(1)} + c_{13}a_{31} = 0$. Затем первое уравнение системы (5.82) заменяют линейной комбинацией первого и третьего уравнений с коэффициентами c_{13} и s_{13} , а третье уравнение — аналогичной комбинацией с коэффициентами $-s_{13}$ и c_{13} . Это преобразование системы эквивалентно умножению слева на матрицу

$$\mathbf{T}_{13} = \begin{bmatrix} c_{13} & 0 & s_{13} & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 \\ -s_{13} & 0 & c_{13} & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & \dots & 0 \\ \dots & & & & & \\ 0 & 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

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

Таким же образом x_1 исключают из уравнений с номерами $i = 4, \dots, m$. В результате 1-го шага (состоящего, как мы видели, из $m - 1$ «малых» шагов) система приводится к виду

$$\begin{aligned} a_{11}^{(m-1)}x_1 + a_{12}^{(m-1)}x_2 + a_{13}^{(m-1)}x_3 + \dots + a_{1m}^{(m-1)}x_m &= b_1^{(m-1)}, \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2m}^{(1)}x_m &= b_2^{(1)}, \\ a_{32}^{(1)}x_2 + a_{33}^{(1)}x_3 + \dots + a_{3m}^{(1)}x_m &= b_3^{(1)}, \\ \dots & \\ a_{m2}^{(1)}x_2 + a_{m3}^{(1)}x_3 + \dots + a_{mm}^{(1)}x_m &= b_m^{(1)}. \end{aligned} \quad (5.85)$$

В матричной записи получаем:

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

где $\mathbf{A}^{(1)} = \mathbf{T}_{1m} \dots \mathbf{T}_{13}\mathbf{T}_{12}\mathbf{A}$, $\mathbf{b}^{(1)} = \mathbf{T}_{1m} \dots \mathbf{T}_{13}\mathbf{T}_{12}\mathbf{b}$.

Здесь и далее через \mathbf{T}_{kl} обозначена матрица элементарного преобразования, отличающаяся от единичной матрицы \mathbf{E} только четырьмя элементами. В ней элементы с индексами (k, k) и (l, l) равны c_{kl} , элемент с индексами (k, l) равен s_{kl} , а элемент с индексами (l, k) равен $-s_{kl}$, причем выполнено условие

$$c_{kl}^2 + s_{kl}^2 = 1. \quad (5.86)$$

Действие матрицы \mathbf{T}_{kl} на вектор \mathbf{x} эквивалентно его повороту вокруг оси, перпендикулярной плоскости Ox_kx_l на угол φ_{kl} такой, что $c_{kl} = \cos \varphi_{kl}$, $s_{kl} = \sin \varphi_{kl}$ (существование такого угла гарантируется равенством (5.86)). Эта геометрическая интерпретация и дала название методу вращений. Операцию умножения на матрицу \mathbf{T}_{kl} часто назы-

вают **плоским вращением** (или **преобразованием Гивенса**). Заметим, что $T_{kl}^T = T_{kl}^{-1}$ и, следовательно, матрица T_{kl} ортогональная.

На 2-м шаге метода вращений, состоящем из $m - 2$ «малых» шагов, из уравнений системы (5.85) с номерами $i = 3, 4, \dots, m$ исключают неизвестное x_2 . Для этого каждое i -е уравнение комбинируют со вторым уравнением. В результате приходим к системе

$$\begin{aligned} a_{11}^{(m-1)}x_1 + a_{12}^{(m-1)}x_2 + a_{13}^{(m-1)}x_3 + \dots + a_{1m}^{(m-1)}x_m &= b_1^{(m-1)}, \\ a_{22}^{(m-1)}x_2 + a_{23}^{(m-1)}x_3 + \dots + a_{2m}^{(m-1)}x_m &= b_2^{(m-1)}, \\ a_{33}^{(2)}x_3 + \dots + a_{3m}^{(2)}x_m &= b_3^{(2)}, \\ \dots \\ a_{m3}^{(2)}x_3 + \dots + a_{mm}^{(2)}x_m &= b_m^{(2)}. \end{aligned}$$

В матричной форме записи получаем

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

где $A^{(2)} = T_{2m} \dots T_{24} T_{23} A^{(1)}$, $\mathbf{b}^{(2)} = T_{2m} \dots T_{24} T_{23} \mathbf{b}^{(1)}$.

После завершения $(m - 1)$ -го шага система принимает вид

$$\begin{aligned} a_{11}^{(m-1)}x_1 + a_{12}^{(m-1)}x_2 + a_{13}^{(m-1)}x_3 + \dots + a_{1m}^{(m-1)}x_m &= b_1^{(m-1)}, \\ a_{22}^{(m-1)}x_2 + a_{23}^{(m-1)}x_3 + \dots + a_{2m}^{(m-1)}x_m &= b_2^{(m-1)}, \\ a_{33}^{(m-1)}x_3 + \dots + a_{3m}^{(m-1)}x_m &= b_3^{(m-1)}, \\ \dots \\ a_{mm}^{(m-1)}x_m &= b_m^{(m-1)}, \end{aligned}$$

или в матричной форме записи

$$A^{(m-1)}\mathbf{x} = \mathbf{b}^{(m-1)},$$

где $A^{(m-1)} = T_{m-1, m} A^{(m-2)}$, $\mathbf{b}^{(m-1)} = T_{m-1, m} \mathbf{b}^{(m-2)}$.

Введем обозначение R для полученной верхней треугольной матрицы $A^{(m-1)}$. Она связана с исходной матрицей A равенством

$$R = TA,$$

где $T = T_{m-1, m} \dots T_{2m} \dots T_{23} T_{1m} \dots T_{13} T_{12}$ — матрица результирующего вращения. Заметим, что матрица T ортогональна как произведение ортогональных матриц. Обозначая $Q = T^{-1} = T^T$, получаем QR-разложение матрицы A .

Обратный ход метода вращений проводится точно так же, как и для метода Гаусса.

Метод вращений обладает замечательной численной устойчивостью. Для него оценка (5.43) справедлива с коэффициентом $f(m) = 6m$. Однако этот метод существенно более трудоемок по сравнению с методом Гаусса. Получение QR -разложения для квадратной матрицы A общего вида требует примерно $2m^3$ арифметических операций.

Пример 5.18. Используя метод вращений, решим на 6-разрядном десятичном компьютере систему уравнений

$$\begin{aligned} 2x_1 - & \quad 9x_2 + 5x_3 = -4, \\ 1.2x_1 - 5.3999x_2 + & \quad 6x_3 = 0.6001, \\ x_1 - & \quad x_2 - 7.5x_3 = -8.5. \end{aligned}$$

Прямой ход.

1-й шаг. Исключим x_1 из второго уравнения. Для этого вычислим c_{12} и s_{12} по формулам (5.80):

$$c_{12} = \frac{2}{\sqrt{2^2 + 1.2^2}} \approx 0.857493, \quad s_{12} = \frac{1.2}{\sqrt{2^2 + 1.2^2}} \approx 0.514495.$$

Преобразуя коэффициенты первого и второго уравнений по формулам (5.83), приходим к системе

$$\begin{aligned} 2.33238x_1 - & \quad 10.4957x_2 + 7.37444x_3 = -3.12122, \\ 7.85493 \cdot 10^{-5}x_2 + & \quad 2.57248x_3 = 2.57256, \\ x_1 - & \quad x_2 - 7.5x_3 = -8.5. \end{aligned}$$

Далее вычислим коэффициенты c_{13} и s_{13} по формулам (5.84):

$$c_{13} = \frac{2.33238}{\sqrt{2.33238^2 + 1^2}} \approx 0.919087, \quad s_{13} = \frac{1}{\sqrt{2.33238^2 + 1^2}} \approx 0.394055.$$

Заменяя первое и третье уравнения их линейными комбинациями с коэффициентами c_{13} , s_{13} и $-s_{13}$, c_{13} соответственно, получаем систему

$$\begin{aligned} 2.53772x_1 - & \quad 10.0405x_2 + 3.82234x_3 = -6.21814, \\ 7.85493 \cdot 10^{-5}x_2 + & \quad 2.57248x_3 = 2.57256, \\ 3.21680x_2 - 9.79909x_3 = & \quad -6.58231. \end{aligned}$$

2-й шаг. В полученной системе имеем $a_{22}^{(1)} = 7.85493 \cdot 10^{-5}$, $a_{32}^{(1)} = 3.21680$. Поэтому

$$c_{23} = \frac{a_{22}^{(1)}}{\sqrt{(a_{22}^{(1)})^2 + (a_{32}^{(1)})^2}} \approx 2.44185 \cdot 10^{-5}, \quad s_{23} \approx 1.00000.$$

Заменяя второе и третье уравнения системы их линейными комбинациями с коэффициентами c_{23} , s_{23} и $-s_{23}$, c_{23} соответственно, приходим к системе

$$\begin{aligned} 2.53772x_1 - 10.0405x_2 + 3.82234x_3 &= -6.21814, \\ 3.21680x_2 + 9.79903x_3 &= -6.58225, \\ -2.57272x_3 &= -2.57272. \end{aligned}$$

Обратный ход дает последовательно значения $x_3 = 1$, $x_2 = 0.999994$, $x_1 = -1.58579 \cdot 10^{-5}$. ▲

2. Метод отражений. *Матрицами Хаусхолдера* (или *отражений*) называются квадратные матрицы вида

$$V = E - 2ww^T,$$

где w — вектор-столбец в \mathbb{R}^m , имеющий единичную длину: $\|w\|_2 = 1$.

Матрица Хаусхолдера симметрична и ортогональна. Действительно,

$$V^T = (E - 2ww^T)^T = E^T - 2(w^T)^T w^T = E - 2ww^T = V;$$

$$V^T V = VV = (E - 2ww^T)(E - 2ww^T) = E - 4ww^T + 4ww^Tww^T = E.$$

Здесь учтено, что $w^Tw = \|w\|_2^2 = 1$.

Умножение на матрицу V называют *преобразованием Хаусхолдера* (или *отражением*). Действие V на вектор x можно интерпретировать как ортогональное отражение вектора в \mathbb{R}^m относительно гиперплоскости, проходящей через начало координат и имеющей нормальный вектор, равный w .

Как и вращения, отражения используются для обращения в нуль элементов преобразуемой матрицы. Однако здесь с помощью одного отражения можно обратить в нуль уже не один элемент матрицы, а целую группу элементов некоторого столбца или строки. Поэтому, являясь почти столь же устойчивым, как и метод вращений, метод отражений позволяет получить QR-разложение квадратной матрицы общего вида примерно за $(4/3)m^3$ арифметических операций, т.е. в полтора раза быстрее. Изложение самого метода можно найти, например, в [9].

Поясним, как получается QR -разложение матрицы A с помощью преобразований Хаусхолдера. Пусть $\mathbf{a} = (a_1, a_2, \dots, a_m)^T$ — произвольный вектор, у которого одна координат a_2, \dots, a_m отлична от нуля. Покажем, что вектор w в преобразовании Хаусхолдера может быть выбран так, чтобы обратились в нуль все координаты вектора $V\mathbf{a}$ кроме первой: $V\mathbf{a} = c\mathbf{e}_1 = c(1, 0, \dots, 0)^T$. Поскольку ортогональное преобразование не меняет евклидову норму вектора, то $c = \|\mathbf{a}\|_2$ и искомое преобразование таково, что

$$V\mathbf{a} = (E - 2ww^T)\mathbf{a} = \mathbf{a} - 2(\mathbf{w}, \mathbf{a})\mathbf{w} = \mathbf{a} - \alpha\mathbf{w} = \|\mathbf{a}\|_2\mathbf{e}_1,$$

где $\alpha = 2(\mathbf{w}, \mathbf{a})$. Таким образом, вектор w следует выбрать так, чтобы

$$\alpha w_1 - a_1 = \pm \|\mathbf{a}\|_2, \quad \alpha w_2 = a_2, \dots, \alpha w_m = a_m,$$

что эквивалентно равенству $w = \alpha^{-1}(\mathbf{a} \pm \|\mathbf{a}\|_2\mathbf{e}_1)$. Вспоминая, что $\alpha = 2(\mathbf{w}, \mathbf{a}) = 2(w_1a_1 + w_2a_2 + \dots + w_ma_m)$, имеем

$$\alpha^2 = 2(a_1 \pm \|\mathbf{a}\|_2)a_1 + 2(a_2^2 + \dots + a_m^2) = \|\mathbf{a} \pm \|\mathbf{a}\|_2\mathbf{e}_1\|_2^2.$$

Таким образом $w = \frac{v}{\|v\|_2}$, где $v = \mathbf{a} \pm \|\mathbf{a}\|_2\mathbf{e}_1$.

Взяв $\mathbf{a} = \mathbf{a}_1$, где \mathbf{a}_1 — первый из столбцов матрицы A , и положив $P_1 = V$, получим

$$A^{(1)} = P_1 A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \dots & a_{1m}^{(1)} \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & \dots & a_{2m}^{(1)} \\ 0 & a_{32}^{(1)} & a_{33}^{(1)} & \dots & a_{3m}^{(1)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & a_{m2}^{(1)} & a_{m3}^{(1)} & \dots & a_{mm}^{(1)} \end{bmatrix}.$$

Далее, взяв вектор $\mathbf{a}_2 = (a_{22}^{(1)}, \dots, a_{m2}^{(1)})^T \in \mathbb{R}^{m-1}$, с помощью преобразования Хаусхолдера V_{m-1} в пространстве $(m-1)$ -мерных векторов можно обнулить все координаты вектора $V_{m-1}\mathbf{a}_2$ кроме первой. Положив

$$P_2 = \begin{pmatrix} 1 & 0 \\ 0 & V_{m-1} \end{pmatrix},$$

получим:

$$A^{(2)} = P_2 A^{(1)} = P_2 P_1 A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \dots & a_{1m}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \dots & a_{2m}^{(2)} \\ 0 & 0 & a_{23}^{(2)} & \dots & a_{2m}^{(2)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & a_{m3}^{(2)} & \dots & a_{mm}^{(2)} \end{bmatrix}.$$

Заметим, что первый столбец и первая строка матрицы при этом преобразовании не меняется.

Выполнив $m - 1$ шаг этого метода, приходим к верхней треугольной матрице

$$A^{(m-1)} = P_{m-1} \dots P_2 P_1 A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \dots & a_{1m}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \dots & a_{2m}^{(2)} \\ 0 & 0 & a_{33}^{(3)} & \dots & a_{3m}^{(3)} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & a_{mm}^{(m-1)} \end{bmatrix}.$$

Поскольку матрицы P_1, \dots, P_{m-1} симметричны и ортогональны, получено разложение

$$A = QR,$$

где матрица $Q = P_1 P_2 \dots P_{m-1}$ — ортогональная, а матрица $R = A^{(m-1)}$ — верхняя треугольная.

Пример 5.19. Применим метод отражений для решения системы уравнений из примера 5.18

$$2x_1 - 9x_2 + 5x_3 = -4,$$

$$1.2x_1 - 5.3999x_2 + 6x_3 = 0.6001,$$

$$x_1 - x_2 - 7.5x_3 = -8.5.$$

Прямой ход. Возьмем первый столбец матрицы системы $a_1 = (2, 1.2, 1)^T$, положим $v = a_1 + \|a_1\|_2 e_1 = (2 + \sqrt{6.44}, 1.2, 1)^T \approx (4.53772, 1.2, 1)^T$ и выберем

$$w = v / \|v\|_2 \approx (0.945545, 0.250049, 0.208375)^T.$$

Построим матрицу Хаусхолдера

$$V = E - 2ww^T =$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - 2 \begin{bmatrix} 0.945545 \\ 0.250049 \\ 0.208375 \end{bmatrix} \begin{bmatrix} 0.945545 & 0.250049 & 0.208375 \end{bmatrix} \approx$$

$$\approx \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - 2 \begin{bmatrix} 0.894055 & 0.236433 & 0.197028 \\ 0.236433 & 0.0625245 & 0.0521040 \\ 0.197028 & 0.0521040 & 0.0434201 \end{bmatrix} \approx$$

$$\approx \begin{bmatrix} -0.788110 & -0.472866 & -0.394056 \\ -0.472866 & 0.874951 & -0.104208 \\ -0.394056 & -0.104208 & 0.913160 \end{bmatrix}.$$

Применим к левой и правой частям системы преобразование Хаусхолдера и получим эквивалентную систему

$$\begin{aligned} -2.53772x_1 + 10.0405x_2 + 3.82233x_3 &= 6.21815, \\ -0.364646x_2 + 3.66694x_3 &= 3.30229, \\ 3.19606x_2 - 9.44423x_3 &= -6.24817. \end{aligned}$$

Возьмем теперь $a_2 = (-0.364646, 3.19606)^T$ и построим двумерное преобразование Хаусхолдера. Здесь уже $v = a_2 + \|a_2\|_2(1, 0)^T = = (-0.364646 + \sqrt{0.364646^2 + 3.19606^2}, 3.19606)^T \approx (2.85215, 3.19606)^T$, $w = v/\|v\|_2 \approx (0.665824, 0.746109)^T$. Двумерная матрица Хаусхолдера имеет вид

$$V_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 2 \begin{bmatrix} 0.665824 \\ 0.746109 \end{bmatrix} \begin{bmatrix} 0.665824 & 0.746109 \end{bmatrix} \approx \begin{bmatrix} 0.113357 & -0.993555 \\ -0.993555 & -0.113357 \end{bmatrix}.$$

Умножив обе части системы на матрицу

$$P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0.113357 & -0.993555 \\ 0 & -0.993555 & -0.113357 \end{bmatrix},$$

получим

$$\begin{aligned} -2.53772x_1 + 10.0405x_2 - 3.82233x_3 &= 6.21815, \\ -3.21680x_2 + 9.79904x_3 &= 6.58224, \\ -2.57274x_3 &= -2.57273. \end{aligned}$$

Обратный ход последовательно дает значения $x_3 \approx 0.999996$, $x_2 \approx 0.999988$, $x_1 \approx -3.35721 \cdot 10^{-5}$. ▲

§ 5.11. Итерационное уточнение

В большинстве случаев метод Гаусса с выбором главного элемента позволяет получить приближенное решение с довольно высокой точностью. Однако иногда возникает необходимость найти решение с большей точностью. Полезно знать, что существует метод, позволяющий найти приближенное решение с относительной точностью, сравнимой с ϵ_m , если только число обусловленности не слишком велико. Этот метод, называемый *итерационным уточнением*, требует небольшого (примерно на 25 %) увеличения машинного времени по сравнению с затратами на получение решения методом Гаусса.

Пусть $\mathbf{x}^{(0)}$ — найденное на компьютере приближенное решение системы $A\mathbf{x} = \mathbf{b}$. Напомним (см. § 5.1), что невязка $\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)}$ и погрешность $\mathbf{e}^{(0)} = \mathbf{x} - \mathbf{x}^{(0)}$ связаны равенством

$$A\mathbf{e}^{(0)} = \mathbf{r}^{(0)}. \quad (5.87)$$

Если бы удалось найти $\mathbf{e}^{(0)}$ как точное решение системы (5.87), то вектор $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \mathbf{e}^{(0)}$ дал бы точное решение системы $A\mathbf{x} = \mathbf{b}$. Однако в действительности вычисленное на компьютере значение $\mathbf{x}^{(1)}$ неизбежно будет содержать погрешность. Тем не менее можно ожидать, что $\mathbf{x}^{(1)}$ окажется лучшим приближением, чем $\mathbf{x}^{(0)}$. Используя приближение $\mathbf{x}^{(1)}$, аналогичным образом можно найти приближение $\mathbf{x}^{(2)}$.

Опишем более подробно очередной $(k+1)$ -й шаг метода.

1°. Вычисляют $\mathbf{r}^{(k)} \approx \mathbf{b} - A\mathbf{x}^{(k)}$. Исключительно важно, чтобы вычисление $\mathbf{r}^{(k)}$ производилось с повышенной точностью. Дело в том, что $\mathbf{b} \approx A\mathbf{x}^{(k)}$ и поэтому при вычислении невязки неизбежно вычитание близких чисел, а следовательно, потеря большого числа значащих цифр. Одна из возможностей состоит в использовании для вычисления $A\mathbf{x}^{(k)}$ и $\mathbf{b} - A\mathbf{x}^{(k)}$ арифметики удвоенной точности.

2° Вычисляют решение системы $A\mathbf{e}^{(k)} = \mathbf{r}^{(k)}$. Так как матрица A не меняется, то получение очередного приближения с использованием

однажды вычисленного LU -разложения матрицы A (с учетом необходимых перестановок строк) требует сравнительно небольшого числа (примерно m^2) арифметических действий.

3°. Вычисляют $x^{(k+1)} \approx x^{(k)} + e^{(k)}$.

Если число обусловленности не очень велико (например, $\text{cond}(A) \ll \varepsilon_m^{-1}$), то метод довольно быстро сходится. Сходимость характеризуется постепенным установлением значащих цифр в приближениях $x^{(k)}$. Если же процесс расходится, то в приближениях не устанавливаются даже старшие значащие цифры.

Замечание. Итерационное уточнение не следует рассматривать как средство, позволяющее справиться с плохой обусловленностью системы. Его применение для решения плохо обусловленных систем целесообразно только, если матрица A и вектор правых частей b в компьютере заданы точно или с повышенной точностью.

Необходимо отметить одну интересную особенность. В процессе итерационного уточнения невязки $r^{(k)}$ обычно не уменьшаются, а даже несколько возрастают по величине.

Кроме того, как оказывается, приближения $x^{(k)}$ могут быть использованы для грубого по порядку оценивания естественного числа обусловленности $v_\delta(x)$. Для этого достаточно воспользоваться приближенной формулой

$$v_\delta(x) \approx \varepsilon_m^{-1} \|e^{(k)}\| / \|x^{(k)}\|. \quad (5.88)$$

Пример 5.20. Используя алгоритм итерационного уточнения, найдем решение системы

$$\begin{aligned} 1.03x_1 + 0.991x_2 &= 2.51, \\ 0.991x_1 + 0.943x_2 &= 2.41 \end{aligned} \quad (5.89)$$

на 3-разрядном десятичном компьютере, считая, что режим вычислений с удвоенной точностью на нем эквивалентен использованию 6-разрядного десятичного компьютера.

Вычислим множитель первого (и в данном случае последнего) шага прямого хода метода Гаусса: $\mu = 0.991/1.03 \approx 0.962$. Так как при вычислении на 3-разрядном компьютере имеем $0.943 - 0.991 \cdot 0.962 \approx -0.01$ и $2.41 - 2.51 \cdot 0.962 \approx 0$, то второе уравнение системы (5.89) приводится к виду $-0.01x_2 = 0$.

Обратная подстановка дает приближенное решение $x_1^{(0)} = 2.44$, $x_2^{(0)} = 0$.

Итерационное уточнение.

1-й шаг. Вычислим (с удвоенной точностью) компоненты невязки:

$$\begin{aligned} r_1^{(0)} &= 2.51 - 1.03x_1^{(0)} - 0.991x_2^{(0)} = -0.00320, \\ r_2^{(0)} &= 2.41 - 0.991x_1^{(0)} - 0.943x_2^{(0)} = -0.00840. \end{aligned}$$

Вычислив методом Гаусса решение системы

$$1.03e_1^{(0)} + 0.991e_2^{(0)} = -0.00320,$$

$$0.991e_1^{(0)} + 0.943e_2^{(0)} = -0.00840,$$

получим $e_1^{(0)} = -0.481$, $e_2^{(0)} = 0.496$. Завершается 1-й шаг вычислением $x_1 = x_1^{(0)} + e_1^{(0)}$, $x_2 = x_2^{(0)} + e_2^{(0)}$, приводящим на 3-разрядном компьютере к значениям $x_1^{(1)} = 1.96$, $x_2^{(1)} = 0.496$.

2-й шаг дает значения $r_1^{(1)} = -0.000336$, $r_2^{(1)} = -0.000088$, $e_1^{(1)} = 0.0223$, $e_2^{(1)} = -0.0235$ и $x_1^{(2)} = 1.98$, $x_2^{(2)} = 0.473$.

3-й шаг дает значения $r_1^{(2)} = 0.00186$, $r_2^{(2)} = 0.00178$, $e_1^{(2)} = 0.00084$, $e_2^{(2)} = 0.001$ и $x_1^{(3)} = 1.98$, $x_2^{(3)} = 0.474$.

Сравнивая $x^{(2)}$ и $x^{(3)}$, замечаем, что последние значения цифры практически установились и, следовательно, процесс следует завершить, приняв $x_1 \approx 1.98$, $x_2 \approx 0.474$.

Использование формулы (5.88) с $\varepsilon_m = 5 \cdot 10^{-4}$ дает следующую оценку естественного числа обусловленности:

$$v_\delta(x) \approx \varepsilon_m^{-1} \|e^{(0)}\|_\infty / \|x^{(0)}\|_\infty \approx 2 \cdot 10^3 \cdot 0.496 / 2.44 \approx 409.$$

Приведем для сравнения действительные значения решения и естественного числа обусловленности: $x_1 \approx 1.9812$, $x_2 \approx 0.4735$, $v_\delta(x) \approx 273$ (см. пример 5.5). ▲

§ 5.12. Линейная задача метода наименьших квадратов

1. Линейная задача метода наименьших квадратов и переопределенные системы. Метод наименьших квадратов широко используется в самых разных предметных областях при обработке экспериментальных данных. В статистике этот метод принято называть линейной регрессией. Инженеры приходят к нему, занимаясь оцениванием параметров или фильтрацией. Мы встретимся с методом наименьших квадратов в гл. 11 при изучении методов приближения таблично заданных функций. Автором метода является Гаусс, который изобрел его в 1795 г. при решении поставленной перед ним немецким правительством задачи обновления координат землемерных вех из государственной базы данных.

Независимо от того, для решения какой прикладной задачи используется метод наименьших квадратов, в плане вычислений он сводится к следующей ключевой задаче, которую принято называть *линейной задачей метода наименьших квадратов*. Пусть A — прямоугольная матрица размера $n \times m$ с элементами a_{ij} , а $\mathbf{b} \in \mathbb{R}^n$ — заданный вектор. Требуется найти вектор $\bar{\mathbf{x}} \in \mathbb{R}^m$, минимизирующий величину $\|A\bar{\mathbf{x}} - \mathbf{b}\|_2$, т.е. такой вектор $\bar{\mathbf{x}}$, что

$$\|A\bar{\mathbf{x}} - \mathbf{b}\|_2 = \min_{\mathbf{x} \in \mathbb{R}^m} \|A\mathbf{x} - \mathbf{b}\|_2. \quad (5.90)$$

Линейная задача метода наименьших квадратов естественным образом возникает и при решении систем линейных алгебраических уравнений. Предположим, что требуется решить систему n линейных алгебраических уравнений с m неизвестными, где $n > m$

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m &= b_1, \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2m}x_m &= b_2, \\ \dots \dots \dots \dots \dots \dots \dots & \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nm}x_m &= b_n. \end{aligned} \quad (5.91)$$

Системы уравнений (5.91), в которых число уравнений превышает число неизвестных, называют *переопределенными*. Для переопределенной системы, вообще говоря, нельзя найти вектор \mathbf{x} , точно удовлетворяющий уравнениям системы. Однако, хотя уравнения системы нельзя удовлетворить точно, можно попытаться удовлетворить их как можно точнее, минимизируя величину вектора невязки $\mathbf{r} = \mathbf{b} - A\mathbf{x}$. Выбор в качестве минимизируемой величины евклидовой нормы невязки $\|\mathbf{r}\|_2 =$

$= \|b - Ax\|_2$ приводит к *методу наименьших квадратов решения переопределенных систем*. Другими словами, в этом методе предлагается за решение переопределенной системы принять решение \bar{x} линейной задачи метода наименьших квадратов, обладающее свойством (5.90).

2. Нормальная система метода наименьших квадратов. Простейшим подходом к решению линейной задачи метода наименьших квадратов является использование *нормальной системы метода наименьших квадратов*

$$A^T A x = A^T b. \quad (5.92)$$

Обоснованием такого подхода является следующий результат.

Теорема 5.5. Вектор \bar{x} минимизирует величину $\|Ax - b\|_2$ тогда и только тогда, когда он является решением нормальной системы (5.92).

Доказательство. Для любых векторов $x, \bar{x} \in \mathbb{R}^m$ имеем

$$Ax - b = A\bar{x} - b + A(x - \bar{x}).$$

Отсюда

$$\begin{aligned} \|Ax - b\|_2^2 &= \|A\bar{x} - b\|_2^2 + 2(A\bar{x} - b, A(x - \bar{x})) + \|A(x - \bar{x})\|_2^2 = \\ &= \|A\bar{x} - b\|_2^2 + 2(z, x - \bar{x}) + \|A(x - \bar{x})\|_2^2, \end{aligned}$$

$$\text{где } z = A^T(A\bar{x} - b) = A^T A \bar{x} - A^T b.$$

Если \bar{x} удовлетворяет нормальной системе (5.92), то $z = 0$ и

$$\|Ax - b\|_2^2 = \|A\bar{x} - b\|_2^2 + \|A(x - \bar{x})\|_2^2 \geq \|A\bar{x} - b\|_2^2$$

для всех $x \in \mathbb{R}^m$. Это означает, что \bar{x} минимизирует $\|Ax - b\|_2$.

Если же \bar{x} не является решением нормальной системы (5.92), то $z \neq 0$. Возьмем $x = \bar{x} - \varepsilon z$ и заметим, что

$$\|Ax - b\|_2^2 = \|A\bar{x} - b\|_2^2 - 2\varepsilon \|z\|_2^2 + \varepsilon^2 \|Az\|_2^2 < \|A\bar{x} - b\|_2^2$$

для всех достаточно малых $\varepsilon > 0$. Это означает, что \bar{x} не обладает свойством (5.90). ■

Обратим внимание на то, что вектор Ax представляет собой линейную комбинацию столбцов a_1, a_2, \dots, a_m матрицы A с коэффициентами x_1, x_2, \dots, x_m :

$$Ax = \sum_{j=1}^m x_j a_j.$$

Таким образом, справедливо следующее утверждение.

Лемма 5.2. Столбцы матрицы A линейно зависимы тогда и только тогда, когда существует ненулевой вектор x такой, что $Ax = 0$. Столбцы матрицы A линейно независимы тогда и только тогда, когда $Ax \neq 0$ для всех $x \neq 0$. ■

Рассмотрим некоторые свойства матрицы $A^T A$ и нормальной системы (5.92). Заметим, что матрица $A^T A$ — квадратная размера $m \times m$.

Лемма 5.3. Матрица $A^T A$ — симметричная, причем

$$(A^T Ax, x) \geq 0 \quad \forall x \in \mathbb{R}^m.$$

Если столбцы матрицы A линейно независимы, то матрица $A^T A$ — положительно определенная.

Доказательство. Матрица симметрична, так как

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

Кроме того,

$$(A^T Ax, x) = (Ax, Ax) = \|Ax\|_2^2 \geq 0 \quad \forall x \in \mathbb{R}^m. \quad (5.93)$$

Если столбцы матрицы A линейно независимы, то в силу леммы 5.2 имеем $Ax \neq 0$ для всех $x \neq 0$. Поэтому

$$(A^T Ax, x) = \|Ax\|_2^2 > 0 \quad \forall x \neq 0, \quad x \in \mathbb{R}^m,$$

что означает положительную определенность матрицы $A^T A$. ■

Лемма 5.4. Вектор $x \in \mathbb{R}^m$ является решением однородной системы

$$A^T Ax = 0 \quad (5.94)$$

тогда и только тогда, когда $Ax = 0$.

Доказательство. Если x является решением системы (5.94), то из равенства (5.93) следует, что $\|Ax\|_2 = 0$, то есть $Ax = 0$.

Если же $Ax = 0$, то умножив обе части этого равенства слева на A^T , придем к (5.94). ■

Лемма 5.5. Матрица $A^T A$ вырождена тогда и только тогда, когда столбцы матрицы A линейно зависимы.

Доказательство. Напомним, что квадратная матрица $A^T A$ вырождена тогда и только тогда, когда однородная система (5.94) имеет нетривиальное решение $x \neq 0$.

Пусть матрица $A^T A$ вырождена и x — нетривиальное решение системы (5.94). Тогда $Ax = 0$ в силу леммы 5.4, и в силу леммы 5.2 столбцы матрицы A линейно зависимы.

Пусть теперь столбцы матрицы A линейно зависимы. Тогда в силу леммы 5.2 существует вектор $x \neq 0$ такой, что $Ax = 0$ и поэтому x является решением системы (5.94). Следовательно, матрица $A^T A$ вырождена. ■

Теорема 5.6. *Нормальная система метода наименьших квадратов (5.94) имеет решение. Это решение единствено тогда и только тогда, когда столбцы матрицы A линейно независимы.*

Доказательство. Пусть $x \in \mathbb{R}^m$ — произвольное решение однородной системы (5.94). Тогда в силу леммы 5.4 $Ax = 0$ и поэтому $(A^T b, x) = (b, Ax) = (b, 0) = 0$.

Таким образом, правая часть нормальной системы (5.92) ортогональна всем решениям однородной системы (5.94). Следовательно, в силу теоремы Фредгольма система (5.92) имеет решение.

Напомним, что решение системы линейных алгебраических уравнений с квадратной матрицей единствено тогда и только тогда, когда матрица системы невырождена. Из леммы 5.5 следует, что решение нормальной системы единствено тогда и только тогда, когда столбцы матрицы A линейно независимы. ■

В силу теоремы 5.5 решения линейной задачи метода наименьших квадратов и нормальной системы (5.94) совпадают. Поэтому теорему 5.6 можно переформулировать следующим образом.

Теорема 5.7. *Линейная задача метода наименьших квадратов имеет решение. Это решение единствено тогда и только тогда, когда столбцы матрицы A линейно независимы.* ■

Как следует из леммы 5.3, в случае, когда столбцы матрицы A линейно независимы, матрица $A^T A$ — симметричная и положительно определенная. Поэтому решение нормальной системы (5.94) можно найти методом Холецкого примерно за $m^3 / 3$ арифметических операций. Правда, для формирования системы (с учетом симметричности матрицы $A^T A$) нужно примерно $n^2 m$ арифметических операций. Если $n \gg m$,

то основные затраты машинного времени будут произведены на этапе формирования матрицы $A^T A$.

Замечание. Если $m \ll n$, то матрица $A^T A$ и вектор $A^T b$ содержат $m^2 \ll n$ и $m \ll n$ элементов. Поэтому процесс их формирования можно рассматривать как сжатие данных.

2. Алгоритм Грама–Шмидта и QR-разложение матрицы A . С представлением квадратной матрицы A в виде произведения ортогональной матрицы Q на верхнюю треугольную матрицу R , то есть в виде

$$A = QR \quad (5.95)$$

мы уже встречались в § 5.10.

В случае, когда A — прямоугольная матрица размера $n \times m$, будем называть *QR-разложением матрицы A* ее представление в виде (5.95), где Q — прямоугольная матрица размера $n \times m$ с ортонормированными столбцами q_1, q_2, \dots, q_m , а R — верхняя треугольная квадратная матрица размера $m \times m$ с элементами r_{ij} .

Замечание. Ортонормированность столбцов матрицы Q означает, что $(q_i, q_j) = 0$ для всех $i \neq j$ и $(q_j, q_j) = 1$ для всех j , что эквивалентно выполнению равенства

$$Q^T Q = E, \quad (5.96)$$

где E — единичная матрица порядка m .

Заметим, что равенство (5.95) означает, что существуют векторы $q_1, q_2, \dots, q_m \in \mathbb{R}^m$ и коэффициенты r_{ik} , $1 \leq i \leq k \leq m$ такие, что для всех $k = 1, 2, \dots, m$ столбцы a_k матрицы A могут быть представлены в виде линейной комбинации

$$a_k = \sum_{i=1}^k r_{ik} q_i.$$

Применим классический *алгоритм Грама¹–Шмидта²* для ортогонализации столбцов a_1, a_2, \dots, a_m матрицы A , предполагая, что они линейно независимы. Напомним, что этот алгоритм состоит из m шагов,

¹ Иорген Педерсен Грам (1850—1916) — датский математик.

² Эрхард Шмидт (1876—1959) — немецкий математик.

на каждом из которых вычисляются очередной вектор \mathbf{q}_k и коэффициенты r_{ik} , $1 \leq i \leq k$.

На первом шаге полагают $r_{11} = \|\mathbf{a}_1\|_2$ и $\mathbf{q}_1 = \mathbf{a}_1/r_{11}$.

Пусть на $(k-1)$ -м шаге найдена ортонормированная система векторов $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{k-1}$. Тогда на k -м шаге производят следующие вычисления:

$$r_{ik} = (\mathbf{a}_k, \mathbf{q}_i), \quad i = 1, 2, \dots, k-1, \quad (5.97)$$

$$\tilde{\mathbf{q}}_k = \mathbf{a}_k - \sum_{i=1}^{k-1} r_{ik} \mathbf{q}_i, \quad (5.98)$$

$$r_{kk} = \|\tilde{\mathbf{q}}_k\|_2, \quad \mathbf{q}_k = \tilde{\mathbf{q}}_k / r_{kk}. \quad (5.99)$$

Заметим, что $\tilde{\mathbf{q}}_k \neq 0$ и $r_{kk} = \|\tilde{\mathbf{q}}_k\|_2 \neq 0$, так как в противном случае

$$\mathbf{a}_k = \sum_{i=1}^{k-1} r_{ik} \mathbf{q}_i. \quad (5.100)$$

Так как каждый из векторов \mathbf{q}_j , $1 \leq j < k$ является линейной комбинацией векторов $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{k-1}$, то равенство (5.100) означает, что вектор \mathbf{a}_k является линейной комбинацией векторов $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{k-1}$, что противоречит предположению о том, что столбцы матрицы A линейно независимы.

Обратим внимание на то, что $\|\mathbf{q}_k\|_2 = 1$ и $(\mathbf{q}_k, \mathbf{q}_j) = 0$ для всех $j = 1, 2, \dots, k-1$. Действительно, в силу формул (5.98), (5.99) имеем

$$r_{kk}(\mathbf{q}_k, \mathbf{q}_j) = (\tilde{\mathbf{q}}_k, \mathbf{q}_j) = (\mathbf{a}_k, \mathbf{q}_j) - \sum_{i=1}^{k-1} r_{ik}(\mathbf{q}_i, \mathbf{q}_j) = r_{jk} - r_{jk} = 0.$$

Таким образом, полученная на k -м шаге алгоритма система векторов $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$ также является ортонормированной.

Теорема 5.8. Пусть A — матрица размера $n \times m$, столбцы которой линейно независимы. Тогда существует единственная матрица Q размера $n \times m$ с ортонормированными столбцами и верхняя треугольная матрица R с положительными диагональными элементами r_{ii} ($1 \leq i \leq m$) такие, что справедливо представление (5.95).

Доказательство. Воспользуемся алгоритмом Грама—Шмидта (5.97)–(5.99) и получим матрицу \mathbf{Q} с ортонормированными столбцами $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_m$ и верхнюю треугольную матрицу \mathbf{R} с элементами r_{ik} , $1 \leq i \leq k \leq m$ (при $k < i$ полагаем $r_{ik} = 0$).

Запишем формулу (5.98) в виде

$$\mathbf{a}_k = \tilde{\mathbf{q}}_k + \sum_{i=1}^{k-1} r_{ik} \mathbf{q}_i = \sum_{i=1}^k r_{ik} \mathbf{q}_i.$$

Последнее равенство означает, что алгоритм Грама—Шмидта дает QR -разложение матрицы \mathbf{A} .

Докажем теперь единственность матриц \mathbf{Q} и \mathbf{R} . Используя разложение (5.95) и равенство (5.96), имеем

$$\mathbf{A}^T \mathbf{A} = [\mathbf{Q} \mathbf{R}]^T \mathbf{Q} \mathbf{R} = \mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R} = \mathbf{R}^T \mathbf{R}.$$

Заметим, что матрица $\mathbf{L} = \mathbf{R}^T$ — нижняя треугольная. Поэтому для матрицы $\mathbf{A}^T \mathbf{A}$ мы получили разложение Холецкого (см. § 5.8):

$$\mathbf{A}^T \mathbf{A} = \mathbf{R}^T \mathbf{R} = \mathbf{L} \mathbf{L}^T. \quad (5.101)$$

Как следует из результатов § 5.8, матрица \mathbf{L} в этом разложении определяется единственным образом. Поэтому матрицы $\mathbf{R} = \mathbf{L}^T$ и $\mathbf{Q} = \mathbf{A} \mathbf{R}^{-1}$ также определяются однозначно. ■

3. Модифицированный алгоритм Грама—Шмидта. К сожалению, классический алгоритм Грама—Шмидта обладает чрезмерно высокой чувствительностью к вычислительной погрешности. Поэтому для численной ортогонализации столбцов матрицы \mathbf{A} используется *модифицированный алгоритм Грама—Шмидта*, устойчивый к вычислительной погрешности.

При отсутствии вычислительной погрешности этот алгоритм дает те же матрицы \mathbf{Q} и \mathbf{R} , что и алгоритм (5.97)–(5.99).

На k -м шаге этого алгоритма получают очередной столбец \mathbf{q}_k матрицы \mathbf{Q} , модифицированные столбцы $\mathbf{a}_{k+1}^{(k)}, \dots, \mathbf{a}_m^{(k)}$ матрицы \mathbf{A} и коэффициенты r_{kj} для $j = k, k+1, \dots, m$.

На первом шаге полагают

$$r_{11} = \|\mathbf{a}_1\|_2, \quad \mathbf{q}_1 = \mathbf{a}_1 / r_{11},$$

$$\mathbf{a}_j^{(1)} = \mathbf{a}_j, \quad j = 2, \dots, m.$$

Пусть на $(k - 1)$ -м шаге получены столбцы $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{k-1}$ и модифицированные столбцы $\mathbf{a}_k^{(k-1)}, \mathbf{a}_{k+1}^{(k-1)}, \dots, \mathbf{a}_m^{(k-1)}$. Тогда на k -м шаге производятся следующие вычисления:

$$\begin{aligned} r_{kk} &= \|\mathbf{a}_k^{(k-1)}\|_2, \quad \mathbf{q}_k = \mathbf{a}_k^{(k-1)} / r_{kk}, \\ r_{kj} &= (\mathbf{a}_j^{(k)}, \mathbf{q}_k), \quad j = k + 1, \dots, m, \\ \mathbf{a}_j^{(k)} &= \mathbf{a}_j^{(k-1)} - r_{kj} \mathbf{q}_k, \quad j = k + 1, \dots, m. \end{aligned}$$

Замечание. Для получения QR -разложения матрицы A могут быть применены и другие методы. Часто используется метод отражений, основанный на преобразованиях Хаусхолдера (см. § 5.10) и обладающий лучшими, чем модифицированный алгоритм Грама—Шмидта свойствами вычислительной устойчивости.

4. Использование QR -разложения для решения линейной задачи метода наименьших квадратов. Пусть QR -разложение матрицы A найдено. Тогда согласно теореме 5.5 вектор \bar{x} , являющийся решением линейной задачи метода наименьших квадратов, удовлетворяет нормальной системе (5.92). Используя разложения (5.95) и (5.101), запишем нормальную систему в виде

$$\mathbf{R}^T \mathbf{R} \mathbf{x} = \mathbf{R}^T \mathbf{Q}^T \mathbf{b}.$$

Положив $\mathbf{c} = \mathbf{Q}^T \mathbf{b}$ и умножив левую и правую части этой системы на $(\mathbf{R}^T)^{-1}$, приходим к системе линейных алгебраических уравнений с верхней треугольной матрицей

$$\mathbf{R} \mathbf{x} = \mathbf{c}, \tag{5.102}$$

решение которой совпадает с решением линейной задачи метода наименьших квадратов.

Таким образом, для решения линейной задачи метода наименьших квадратов можно сначала модифицированным методом Грама—Шмидта найти QR -разложение матрицы A , затем вычислить вектор $\mathbf{c} = \mathbf{Q}^T \mathbf{b}$ и, наконец, найти \bar{x} как решение системы (5.102).

Однако, если формально действовать по этой схеме, то на этапе вычисления вектора \mathbf{c} умножением матрицы \mathbf{Q}^T на вектор \mathbf{b} можно потерять все преимущества в вычислительной устойчивости, которые были получены переходом от классического алгоритма Грама—Шмидта к его модифицированному варианту.

Для сохранения хороших свойств вычислительной устойчивости вектор $\mathbf{c} = (c_1, c_2, \dots, c_m)^T$ вычисляют иначе. Сначала полагают $\mathbf{r}^{(0)} = \mathbf{b}$, а затем для $k = 1, 2, \dots, m$ проводят следующие вычисления

$$c_k = (\mathbf{r}^{(k-1)}, \mathbf{q}_k), \quad \mathbf{r}^{(k)} = \mathbf{r}^{(k-1)} - c_k \mathbf{q}_k.$$

В результате получают вектор \mathbf{c} и невязку $\mathbf{r} = \mathbf{r}^{(m)} = \mathbf{b} - \mathbf{A}\bar{\mathbf{x}}$.

Обратим внимание на то, что основные вычислительные затраты при решении линейной задачи метода наименьших квадратов производятся на этапе получения QR -разложения. В итоге время, необходимое для решения задачи с помощью QR -разложения примерно вдвое превышает время, необходимое для решения с использованием нормальной системы.

5. Сингулярное разложение матрицы. Пусть \mathbf{A} — вещественная матрица размера $n \times m$, где $n \geq m$. Тогда существуют такие ортогональные матрицы U и V размера $n \times n$ и $m \times m$ соответственно, что

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T. \quad (5.103)$$

Здесь

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \dots & \sigma_m \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

диагональная матрица размера $n \times m$ с элементами $\sigma_1 \geq \sigma_2 \geq \sigma_3 \dots \geq \sigma_m \geq 0$ на главной диагонали.

Разложение (5.103) называется *сингулярным разложением матрицы A* или *SVD-разложением*¹. Числа $\sigma_1, \dots, \sigma_m$ называются *сингулярными числами* матрицы \mathbf{A} . Эти числа определяются по матрице \mathbf{A} однозначно и являются важными ее характеристиками. Например, число ненулевых сингулярных чисел совпадает с рангом матрицы.

¹ От английского *singular value decomposition*.

Сингулярное разложение существует более 100 лет. Оно было независимо открыто в 1873 г. Бельтрами¹ и в 1874 г. Жорданом² для квадратных матриц. В 30-е годы XX в. оно было распространено на прямоугольные матрицы. В вычислительную практику SVD вошло в 60-е годы и стало важнейшим инструментом вычислений.

6. Использование сингулярного разложения для решения линейной задачи метода наименьших квадратов. Пусть разложение (5.103) получено. Пусть ранг матрицы A равен r , т.е. пусть $\sigma_1 \geq \dots \geq \sigma_r > 0$ и $\sigma_{r+1} = \dots = \sigma_m = 0$. Нашей целью является минимизация величины $\|r\|_2 = \|b - Ax\|_2$.

Так как матрица U ортогональна, то норма невязки r не изменится после ее умножения слева на U^T . Таким образом,

$$\|Ax - b\|_2^2 = \|U^T(Ax - b)\|_2^2 = \|U^T U \Sigma V^T x - U^T b\|_2^2 = \|\Sigma V^T x - U^T b\|_2^2.$$

Обозначая $y = V^T x$ и $d = U^T b$, мы приходим к эквивалентной задаче минимизации величины

$$\begin{aligned}\|\Sigma y - d\|_2^2 &= (\sigma_1 y_1 - d_1)^2 + (\sigma_2 y_2 - d_2)^2 + \dots \\ &\dots + (\sigma_r y_r - d_r)^2 + d_{r+1}^2 + \dots + d_n^2.\end{aligned}$$

Ясно, что минимум равен $d_{r+1}^2 + \dots + d_m^2$ и достигается он при $y_1 = d_1/\sigma_1, \dots, y_r = d_r/\sigma_r$. После того как оказался найденным вектор y , осталось положить $x = Vy$. Если $r = m$, то решение найдено однозначно. Если же $r < m$, то компоненты y_{r+1}, \dots, y_m могут быть выбраны произвольным образом и, следовательно, задача имеет бесконечно много решений. Выбор $y_{r+1} = \dots = y_m = 0$ дает решение x с минимальной нормой (это решение называется *псевдорешением* системы (5.91)).

Понятно, что наличие малых сингулярных чисел σ_i приводит к тому, что малые погрешности в задании величин d , приводят к большим погрешностям в вычислении компонент $y_i = d_i/\sigma_i$. В этом случае задача является плохо обусловленной. Правильное использование SVD предполагает наличие некоторого допуска $\sigma^* > 0$. Все сингулярные числа

¹ Эудженио Бельтрами (1835—1900) — итальянский математик.

² Камиль Мари Эдмон Жордан (1838—1922) — французский математик, внесший значительный вклад в алгебру, теорию чисел, теорию функций, геометрию, топологию и дифференциальные уравнения.

$\sigma_i \leq \sigma_*$ считаются пренебрежимо малыми и заменяются нулями, а соответствующие значения y_i полагаются равными нулю.

Заметим, что увеличение значения параметра σ_* приводит к увеличению нормы невязки; однако при этом уменьшается норма решения и оно становится менее чувствительным к входным данным.

7. Дополнительная информация о сингулярном разложении. Обозначим столбцы матрицы U через u_1, \dots, u_n , а столбцы матрицы V — через v_1, \dots, v_m . Векторы u_i и v_i называются соответственно *левыми* и *правыми сингулярными векторами* матрицы A .

Равенство (5.90) можно записать в виде $AV = U\Sigma$ или $A^T U = V\Sigma^T$. Сравнивая эти матричные равенства по столбцам, видим, что сингулярные векторы обладают следующими свойствами:

$$A v_i = \sigma_i u_i, \quad A^T u_i = \sigma_i v_i, \quad 1 \leq i \leq m.$$

Из этих равенств следует, в частности, что

$$A^T A v_i = \sigma_i^2 v_i, \quad 1 \leq i \leq m.$$

Таким образом, квадраты сингулярных чисел матрицы A совпадают с собственными значениями матрицы $A^T A$. Для симметричных матриц сингулярные числа просто равны модулям собственных значений.

Обратим теперь внимание на норму матрицы A . Делая замену $y = Vx$ и учитывая, что умножение вектора на ортогональную матрицу не меняет его евклидовую норму, имеем:

$$\begin{aligned} \|A\|_2 &= \max_{\|x\|_2=1} \|Ax\|_2 = \max_{\|x\|_2=1} \|U^T Ax\|_2 = \max_{\|x\|_2=1} \|\Sigma V^T x\|_2 = \\ &= \max_{\|y\|_2=1} \|\Sigma y\|_2 = \max_{\|y\|_2=1} \sqrt{\sigma_1^2 y_1^2 + \dots + \sigma_m^2 y_m^2} = \sigma_1. \end{aligned}$$

Отметим также, что отношение σ_1/σ_m часто используется в качестве числа обусловленности матрицы A . Для квадратной невырожденной матрицы A действительно $\text{cond}_2(A) = \sigma_1/\sigma_m$.

Более подробно о сингулярном разложении и его разнообразных приложениях см. в [31, 54].

8. Сравнение различных методов решения линейной задачи метода наименьших квадратов. Мы рассмотрели следующие методы решения линейной задачи метода наименьших квадратов:

- 1) использование нормальной системы метода наименьших квадратов;

- 2) использование QR -разложения;
- 3) использование сингулярного разложения матрицы A .

Какой метод следует предпочесть при решении линейной задачи метода наименьших квадратов? Наиболее быстрым является метод нормальных уравнений. Затем идет метод, основанный на использовании QR -разложения. Наконец, наибольших вычислительных затрат требует использование сингулярного разложения.

Если матрица A хорошо обусловлена, то следует предпочесть метод нормальных уравнений. Однако, если матрица A плохо обусловлена, нормальная система имеет число обусловленности, примерно равное квадрату числа обусловленности матрицы A . Поэтому можно ожидать, что при использовании этого метода будет потеряно вдвое больше значащих цифр по сравнению с методами, основанными на QR -разложении или сингулярном разложении матрицы.

В случае, когда матрица A плохо обусловлена, но система её столбцов линейно независима и не близка к линейно зависимой, наиболее подходящим методом является, по-видимому, использование QR -разложения.

Наиболее медленным, но наиболее надежным методом является использование сингулярного разложения. Особенно важно использование сингулярного разложения в случае, когда система столбцов матрицы A линейно зависима или близка к линейно зависимой.

Таким образом, выбор метода решения зависит от того, что важнее для пользователя: скорость решения задачи или надежность полученных результатов.

§ 5.13. Дополнительные замечания

1. Более подробное и в то же время доступное изложение рассмотренных в этой главе вопросов можно найти, например, в [46, 47, 81, 107, 5*, 16*, 20*]. Особое внимание следует обратить на книгу [31], представляющую собой удачное сочетание учебного пособия и справочника по методам матричных вычислений.

2. Для решения систем линейных алгебраических уравнений создано высококачественное и бесплатно распространяемое программное обеспечение. Пакеты BLAS, LINPACK, EISPACK и LAPACK, ориентированные на системы уравнений с заполненными матрицами, доступны по адресу <http://www.netlib.org>. Пакет SPARSKIT, предназначенный для решения систем уравнений с разреженными матрицами, доступен по адресу <http://www.cs.umn.edu/Research/arpa/SPARSKIT/>.

Обратим также внимание на пакет LinPar, доступный по адресу <http://www.ict.nsc.ru/lipar/>.

3. Интересной модификацией метода Гаусса является *алгоритм Краута*, основанный на изменении порядка выполнения арифметических операций. Существуют два случая, когда эта модификация может иметь преимущества. Во-первых, алгоритм Краута выгоден тогда, когда для вычислений используется калькулятор, так как алгоритм позволяет избегать выписывания промежуточных результатов. Во-вторых, он приводит к наименьшей вычислительной погрешности при условии, что вычисления ведутся на компьютере, обладающем сравнительно быстрой арифметикой с расширенной точностью или особенно быстрой операцией вычисления скалярных произведений $(x, y) = \sum_{i=1}^m x_i y_i$. Обсуждение модификации Краута можно найти, например, в [81, 107].

4. Иногда вместо метода Гаусса используют *метод Гаусса—Жордана*, отличающийся от метода Гаусса тем, что на k -м шаге прямого метода обнуляются не только элементы k -го столбца матрицы, находящиеся под главной диагональю, но и элементы k -го столбца, находящиеся над главной диагональю. В результате матрица системы приводится не к нижнему треугольному виду, а к диагональному виду.

Отметим, что прямой ход метода Гаусса—Жордана требует примерно $m^3/2$ арифметических операций и по затратам машинного времени проигрывает методу Гаусса примерно в полтора раза.

5. Следует отметить, что в данной главе оказались практически не отраженными прямые методы решения систем уравнений с разреженными матрицами (исключение составляет § 5.9, посвященный методу прогонки). Желающим найти доступное изложение современных прямых методов, предназначенных для решения очень больших линейных систем с разреженными матрицами, можно посоветовать обратиться к [48]. Укажем также на книги [38, 77, 116], специально посвященные технологии разреженных матриц.

6. Предположим, что после того, как было найдено решение x системы $Ax = b$, возникла необходимость решить систему $\tilde{A}\tilde{x} = b$ с матрицей \tilde{A} , отличающейся от матрицы A несколькими элементами. Например, могло выясниться, что заданные ранее элементы содержали грубые ошибки. Возможно также, что решаемая задача такова, что в ней элементы матрицы последовательно меняются, а правая часть остается неизменной.

Разумеется, \tilde{x} можно найти, решив систему снова и проигнорировав полученную на предыдущем этапе информацию. Однако в рассматриваемом случае можно найти поправку $\Delta x = \tilde{x} - x$ к найденному ранее решению, использовав всего $O(m^2)$ операции.

Пусть $\tilde{A} = A - uv^T$, где u и v — векторы размерности m . Справедлива формула Шерманна—Моррисона

$$\tilde{A}^{-1} = A^{-1} + \alpha(A^{-1}u)(v^TA^{-1}),$$

где $\alpha = 1/(1 - v^TA^{-1}u)$. Применяя ее, приходим к следующему алгоритму:

- 1°. Систему $Ay = u$ решают относительно y .
- 2°. Систему $A^Tz = v$ решают относительно z .
- 3°. Вычисляют $\alpha = 1/(1 - v^Ty)$, $\beta = z^Tb$ и $\Delta x = \alpha\beta y$.
- 4°. Полагают $\tilde{x} = x + \Delta x$.

Суммарное число операций будет действительно составлять $O(m^2)$, если для решения систем $Ay = u$ и $A^Tz = v$ применить обратную подстановку с использованием LU-разложения матрицы A , найденного ранее на этапе решения системы $Ax = b$.

В случае, когда матрица \tilde{A} отличается от матрицы A только одним элементом $\tilde{a}_{ij} = a_{ij} + \Delta a_{ij}$, можно положить $u = \Delta a_{ij} e_i$, $v = e_j$, а указанный алгоритм упростить следующим образом:

- 1°. Систему $Ay = \Delta a_{ij} e_i$ решают относительно y .
- 2°. Систему $A^Tz = e_j$ решают относительно z .
- 3°. Вычисляют $\alpha = 1/(1 - y_j)$, $\beta = z^Tb$ и $\Delta x = \alpha\beta y$.
- 4°. Полагают $\tilde{x} = x + \Delta x$.

Основываясь на формуле Шерманна—Моррисона, можно указать и способ пересчета LU-разложения матрицы.

О формуле Шерманна—Моррисона и ее применениях см. [54, 31].

Глава 6

ИТЕРАЦИОННЫЕ МЕТОДЫ РЕШЕНИЯ СИСТЕМ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ

Системы линейных алгебраических уравнений можно решать как с помощью прямых, так и итерационных методов. Для систем уравнений средней размерности чаще используют прямые методы.

Итерационные методы применяют главным образом для решения задач большой размерности, когда использование прямых методов невозможно из-за ограничений в доступной оперативной памяти компьютера или из-за необходимости выполнения чрезмерно большого числа арифметических операций. Большие системы уравнений, возникающие в приложениях, как правило, являются разреженными. Методы исключения для решения систем с разреженными матрицами неудобны, например тем, что при их использовании большое число нулевых элементов превращается в ненулевые и матрица теряет свойство разреженности. В противоположность им при использовании итерационных методов в ходе итерационного процесса матрица не меняется, и она, естественно, остается разреженной. Большая эффективность итерационных методов по сравнению с прямыми методами тесно связана с возможностью существенного использования разреженности матриц.

Применение итерационных методов для эффективного решения большой системы уравнений требует серьезного использования ее структуры, специальных знаний и определенного опыта. Именно поэтому разработано большое число различных итерационных методов, каждый из которых ориентирован на решение сравнительно узкого класса задач, и существует довольно мало стандартных программ, реализующих эти методы.

В этой главе будут рассмотрены наиболее простые и известные итерационные методы, позволяющие решать достаточно широкий класс систем.

§ 6.1. Метод простой итерации

1. Приведение системы к виду, удобному для итераций. Для того чтобы применить метод простой итерации к решению системы линейных алгебраических уравнений

$$Ax = b \quad (6.1)$$

с квадратной невырожденной матрицей A , необходимо предварительно преобразовать эту систему к виду

$$\mathbf{x} = \mathbf{Bx} + \mathbf{c}. \quad (6.2)$$

Здесь \mathbf{B} — квадратная матрица с элементами b_{ij} ($i, j = 1, 2, \dots, m$); \mathbf{c} — вектор-столбец с элементами c_i ($i = 1, 2, \dots, m$).

В развернутой форме записи системы (6.2) имеет следующий вид:

$$\begin{aligned} x_1 &= b_{11}x_1 + b_{12}x_2 + b_{13}x_3 + \dots + b_{1m}x_m + c_1, \\ x_2 &= b_{21}x_1 + b_{22}x_2 + b_{23}x_3 + \dots + b_{2m}x_m + c_2, \\ &\dots \\ x_m &= b_{m1}x_1 + b_{m2}x_2 + b_{m3}x_3 + \dots + b_{mm}x_m + c_m. \end{aligned} \quad (6.3)$$

Как правило, операция приведения системы к виду, удобному для итераций (т.е. к виду (6.2)), не является простой и требует специальных знаний, а также существенного использования специфики системы. В некоторых случаях в таком преобразовании нет необходимости, так как сама исходная система уже имеет вид (6.2).¹

Самый простой способ приведения системы к виду, удобному для итераций, состоит в следующем. Из первого уравнения системы (6.1) выразим неизвестное x_1 :

$$x_1 = a_{11}^{-1}(b_1 - a_{12}x_2 - a_{13}x_3 - \dots - a_{1m}x_m),$$

из второго уравнения — неизвестное x_2 :

$$x_2 = a_{22}^{-1}(b_2 - a_{21}x_1 - a_{23}x_3 - \dots - a_{2m}x_m),$$

и т.д. В результате получим систему

$$\begin{aligned} x_1 &= b_{12}x_2 + b_{13}x_3 + \dots + b_{1,m-1}x_{m-1} + b_{1m}x_m + c_1, \\ x_2 &= b_{21}x_1 + b_{23}x_3 + \dots + b_{2,m-1}x_{m-1} + b_{2m}x_m + c_2, \\ x_3 &= b_{31}x_1 + b_{32}x_2 + \dots + b_{3,m-1}x_{m-1} + b_{3m}x_m + c_3, \\ &\dots \\ x_m &= b_{m1}x_1 + b_{m2}x_2 + b_{m3}x_3 + \dots + b_{m,m-1}x_{m-1} + c_m, \end{aligned} \quad (6.4)$$

в которой на главной диагонали матрицы \mathbf{B} находятся нулевые элементы. Остальные элементы вычисляются по формулам

$$b_{ij} = -a_{ij}/a_{ii}, \quad c_i = b_i/a_{ii} \quad (i, j = 1, 2, \dots, m, j \neq i). \quad (6.5)$$

¹ Такие системы возникают, например, при аппроксимации интегральных уравнений Фредгольма второго рода (см. гл. 16).

Конечно, для возможности выполнения указанного преобразования необходимо, чтобы диагональные элементы матрицы A были ненулевыми.

Часто систему (6.1) преобразуют к виду $\mathbf{x} = \mathbf{x} - \tau(\mathbf{A}\mathbf{x} - \mathbf{b})$, где τ — специально выбираемый числовой параметр (см. § 6.4).

2. Описание метода простой итерации. Выберем начальное приближение $\mathbf{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_m^{(0)})$. Подставляя его в правую часть системы (6.2) и вычисляя полученное выражение, находим первое приближение

$$\mathbf{x}^{(1)} = \mathbf{B}\mathbf{x}^{(0)} + \mathbf{c}.$$

Подставляя приближение $\mathbf{x}^{(1)}$ в правую часть системы (6.2), получаем

$$\mathbf{x}^{(2)} = \mathbf{B}\mathbf{x}^{(1)} + \mathbf{c}.$$

Продолжая этот процесс далее, получаем последовательность $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}, \dots$ приближений, вычисляемых по формуле

$$\mathbf{x}^{(k+1)} = \mathbf{B}\mathbf{x}^{(k)} + \mathbf{c}, \quad k = 0, 1, 2, \dots \quad (6.6)$$

В развернутой форме записи формула (6.6) выглядит так:

$$\begin{aligned} x_1^{(k+1)} &= b_{11}x_1^{(k)} + b_{12}x_2^{(k)} + b_{13}x_3^{(k)} + \dots + b_{1m}x_m^{(k)} + c_1, \\ x_2^{(k+1)} &= b_{21}x_1^{(k)} + b_{22}x_2^{(k)} + b_{23}x_3^{(k)} + \dots + b_{2m}x_m^{(k)} + c_2, \\ &\dots \\ x_m^{(k+1)} &= b_{m1}x_1^{(k)} + b_{m2}x_2^{(k)} + b_{m3}x_3^{(k)} + \dots + b_{mm}x_m^{(k)} + c_m. \end{aligned} \quad (6.7)$$

Когда для итераций используется система (6.4) с коэффициентами, вычисленными по формулам (6.5), метод простой итерации принято называть *методом Якоби*¹.

3. Сходимость метода простой итерации

Теорема 6.1. Пусть выполнено условие

$$\|\mathbf{B}\| < 1. \quad (6.8)$$

Тогда:

1) решение $\bar{\mathbf{x}}$, системы (6.2) существует и единствено;

¹ Карл Густав Якоб Якоби (1804—1851) — немецкий математик.

2) при произвольном начальном приближении $\mathbf{x}^{(0)}$ метод простой итерации сходится и справедлива оценка погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| \leq \|\mathbf{B}\|^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|. \quad (6.9)$$

Доказательство. 1) Из курса линейной алгебры известно, что система линейных алгебраических уравнений имеет единственное решение при любой правой части тогда и только тогда, когда соответствующая однородная система имеет только нулевое решение. Пусть \mathbf{x} — решение однородной системы $\mathbf{x} = \mathbf{Bx}$. Тогда $\|\mathbf{x}\| = \|\mathbf{Bx}\| \leq \|\mathbf{B}\| \cdot \|\mathbf{x}\|$. Так как по условию $\|\mathbf{B}\| < 1$, это неравенство возможно только при $\|\mathbf{x}\| = 0$. Следовательно, $\mathbf{x} = 0$ и тем самым первое утверждение теоремы доказано.

2) Вычитая из равенства (6.6) равенство $\bar{\mathbf{x}} = \mathbf{B}\bar{\mathbf{x}} + \mathbf{c}$, получаем

$$\mathbf{x}^{(k+1)} - \bar{\mathbf{x}} = \mathbf{B}(\mathbf{x}^{(k)} - \bar{\mathbf{x}}). \quad (6.10)$$

Вычисляя норму левой и правой частей этого равенства и используя неравенство $\|\mathbf{B}(\mathbf{x}^{(k)} - \bar{\mathbf{x}})\| \leq \|\mathbf{B}\| \|\mathbf{x}^{(k)} - \bar{\mathbf{x}}\|$, имеем $\|\mathbf{x}^{(k+1)} - \bar{\mathbf{x}}\| \leq \|\mathbf{B}\| \|\mathbf{x}^{(k)} - \bar{\mathbf{x}}\|$. Так как это неравенство верно для всех $k \geq 0$, то

$$\begin{aligned} \|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| &\leq \|\mathbf{B}\| \|\mathbf{x}^{(n-1)} - \bar{\mathbf{x}}\| \leq \|\mathbf{B}\|^2 \|\mathbf{x}^{(n-2)} - \bar{\mathbf{x}}\| \leq \dots \\ &\dots \leq \|\mathbf{B}\|^{n-1} \|\mathbf{x}^{(1)} - \bar{\mathbf{x}}\| \leq \|\mathbf{B}\|^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|. \end{aligned}$$

Справедливость неравенства (6.9) установлена. Учитывая, что $\|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|$ не зависит от n , а $\|\mathbf{B}\|^n \rightarrow 0$ при $n \rightarrow \infty$, получаем из него, что $\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| \rightarrow 0$ при $n \rightarrow \infty$. ■

Замечание 1. Теорема 6.1 дает простое достаточное условие (6.8) сходимости метода простой итерации. Грубо это условие можно интерпретировать как условие достаточной малости элементов b_{ij} матрицы \mathbf{B} в системе, приведенной к виду (6.2).

Замечание 2. Если $\|\mathbf{B}\| = \|\mathbf{B}\|_\infty$, то условие (6.8) принимает вид

$$\|\mathbf{B}\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^m |b_{ij}| < 1.$$

Для метода Якоби это условие в силу равенств $b_{ii} = 0$, $b_{ij} = -a_{ij}/a_{ii}$ для $i \neq j$ эквивалентно условию строчного диагонального преобладания (5.45).

Таким образом, для сходимости метода Якоби достаточно, чтобы матрица \mathbf{A} была близка к диагональной.

Отметим еще одно условие *столбцового диагонального преобразования*

$$\sum_{\substack{i=1 \\ i \neq j}}^m |a_{ij}| < |a_{jj}|, \quad j = 1, 2, \dots, m. \quad (6.11)$$

Правда, оно не является достаточным для сходимости метода Якоби.

Замечание 3. Из оценки (6.9) следует, что при выполнении условия (6.8) метод простой итерации сходится со скоростью геометрической прогрессии, знаменатель которой $q = \|\mathbf{B}\|$. Скорость сходимости тем выше, чем меньше величина $\|\mathbf{B}\|$. Хотя метод сходится при любом начальном приближении $\mathbf{x}^{(0)}$, из оценки (6.9) можно сделать полезный вывод: начальное приближение желательно выбирать близким к решению.

Замечание 4. Оценка погрешности (6.9) является априорной. Ее использование для формулировки критерия окончания итераций затруднительно, так как значение $\|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|$ неизвестно, а его грубое оценивание заведомо приведет к завышению необходимого числа итераций.

4. Апостериорная оценка погрешности

Предложение 6.1. Если выполнено условие (6.8), то справедлива апостериорная оценка погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| \leq \frac{\|\mathbf{B}\|}{1 - \|\mathbf{B}\|} \|\mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}\|. \quad (6.12)$$

Доказательство. Запишем равенство (6.10) при $k = n - 1$ в виде

$$\mathbf{x}^{(n)} - \bar{\mathbf{x}} = \mathbf{B}(\mathbf{x}^{(n-1)} - \mathbf{x}^{(n)}) + \mathbf{B}(\mathbf{x}^{(n)} - \bar{\mathbf{x}}).$$

Тогда

$$\begin{aligned} \|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| &= \|\mathbf{B}(\mathbf{x}^{(n-1)} - \mathbf{x}^{(n)}) + \mathbf{B}(\mathbf{x}^{(n)} - \bar{\mathbf{x}})\| \leq \\ &\leq \|\mathbf{B}\| \|\mathbf{x}^{(n-1)} - \mathbf{x}^{(n)}\| + \|\mathbf{B}\| \|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|. \end{aligned}$$

Для завершения доказательства осталось заметить, что полученное неравенство эквивалентно неравенству (6.12). ■

Замечание. Величину, стоящую в правой части неравенства (6.12), можно легко вычислить после нахождения очередного приближения $\mathbf{x}^{(n)}$. Если требуется найти решение с точностью ε , то в силу (6.12) следует вести итерации до выполнения неравенства

$$\frac{\|\mathbf{B}\|}{1 - \|\mathbf{B}\|} \|\mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}\| < \varepsilon. \quad (6.13)$$

Таким образом, в качестве критерия окончания итерационного процесса может быть использовано неравенство

$$\|\mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}\| < \varepsilon_1, \quad (6.14)$$

где $\varepsilon_1 = \frac{1 - \|B\|}{\|B\|} \varepsilon$.

В практике вычислений иногда используют привлекательный своей простотой критерий окончания

$$\|\mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}\| < \varepsilon. \quad (6.15)$$

Отметим, что для метода простой итерации его применение обосновано только тогда, когда $\|B\| \leq 1/2$ (в этом случае $(1 - \|B\|)/\|B\| \geq 1$ и выполнение неравенства (6.15) влечет за собой выполнение неравенства (6.14)). Однако в большинстве реальных случаев величина $\|B\|$ оказывается близкой к единице и поэтому $(1 - \|B\|)/\|B\| \ll 1$. В этих случаях $\varepsilon_1 \ll \varepsilon$ и использование критерия (6.15) приводит к существенно преждевременному окончанию итераций. Величина $\|\mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}\|$ здесь оказывается малой не потому, что приближения $\mathbf{x}^{(n-1)}$ и $\mathbf{x}^{(n)}$ близки к решению, а потому, что метод сходится медленно (ср. с замечанием на с. 104).

Пример 6.1. Использовав метод простой итерации в форме Якоби, найдем решение системы

$$\begin{aligned} 6.25x_1 - x_2 + 0.5x_3 &= 7.5, \\ -x_1 + 5x_2 + 2.12x_3 &= -8.68, \\ 0.5x_1 + 2.12x_2 + 3.6x_3 &= -0.24, \end{aligned} \quad (6.16)$$

с точностью $\varepsilon = 10^{-3}$ в норме $\|\cdot\|_\infty$.

Вычислив коэффициенты по формулам (6.5), приведем систему к виду (6.4)

$$\begin{aligned} x_1 &= 0.16x_2 - 0.08x_3 + 1.2, \\ x_2 &= 0.2x_1 - 0.424x_3 - 1.736, \\ x_3 &= -0.1389x_1 - 0.5889x_2 - 0.0667. \end{aligned} \quad (6.17)$$

В последнем уравнении коэффициенты даны с точностью до погрешности округления. Здесь

$$B = \begin{bmatrix} 0 & 0.16 & -0.08 \\ 0.2 & 0 & -0.424 \\ -0.1389 & -0.5889 & 0 \end{bmatrix}, \quad c = \begin{bmatrix} 1.2 \\ -1.736 \\ -0.0667 \end{bmatrix}. \quad (6.18)$$

Достаточное условие сходимости метода простой итерации выполнено, так как $\|B\|_\infty = \max\{0.24, 0.624, 0.7278\} = 0.7278 < 1$.

Примем за начальное приближение к решению вектор $x^{(0)} = (0, 0, 0)^T$ и будем вести итерации по формуле (6.6) до выполнения критерия окончания (6.14), где в данном случае $\varepsilon_1 = (1 - 0.7278)/0.7278 \cdot 10^{-3} \approx 0.37 \cdot 10^{-3}$. Значения приближения в табл. 6.1 приводятся с четырьмя цифрами после десятичной точки.

Таблица 6.1

n	0	1	2	3	4
$x_1^{(n)}$	0.0000	1.2000	0.9276	0.9020	0.8449
$x_2^{(n)}$	0.0000	-1.7360	-1.4677	-1.8850	-1.8392
$x_3^{(n)}$	0.0000	-0.0667	0.7890	0.6688	0.9181
$\ x^{(n)} - x^{(n-1)}\ _\infty$	—	1.7360	0.8557	0.4173	0.2493

Продолжение табл. 6.1

n	...	12	13	14	15
$x_1^{(n)}$...	0.8006	0.8003	0.8002	0.8001
$x_2^{(n)}$...	-1.9985	-1.9993	-1.9995	-1.9998
$x_3^{(n)}$...	0.9987	0.9990	0.9995	0.9997
$\ x^{(n)} - x^{(n-1)}\ _\infty$...	0.0018	0.0008	0.0005	0.0003

При $n = 15$ условие (6.14) выполняется и можно положить $x_1 = -0.800 \pm 0.001$, $x_2 = -2.000 \pm 0.001$, $x_3 = 1.000 \pm 0.001$. В данном случае точные значения решения $x_1 = -0.8$, $x_2 = -2$, $x_3 = 1$ нам известны (см. пример 5.16).

Заметим, что в действительности решение с точностью до $\varepsilon = 10^{-3}$ было получено уже при $n = 13$. ▲

5. Необходимое и достаточное условие сходимости. Поскольку условие (6.8) является лишь грубым достаточным условием сходимости метода простой итерации, представляет интерес формулировка необходимого и достаточного условия сходимости.

Пусть $\lambda_1, \lambda_2, \dots, \lambda_m$ — собственные значения матрицы B . Спектральным радиусом матрицы B называется величина $\rho(B) = \max_{1 \leq i \leq m} |\lambda_i|$.

Теорема 6.2. *Метод простой итерации (6.6) сходится при любом начальном приближении $x^{(0)}$ тогда и только тогда, когда $\rho(B) < 1$, т.е. когда все собственные значения матрицы B по модулю меньше единицы.*

Замечание. Для симметричной матрицы B ее спектральный радиус $\rho(B)$ совпадает с $\|B\|_2$. Поэтому для сходимости метода простой итерации с симметричной матрицей B условие $\|B\|_2 < 1$ является необходимым и достаточным.

6. Влияние ошибок округления. Из-за наличия ошибок округления реально вычисляемые на компьютере приближения $\tilde{x}^{(n)}$ отличаются от идеальных приближений $x^{(n)}$. Поэтому нельзя утверждать, что для любого $\varepsilon > 0$ найдется номер $n_0(\varepsilon)$, начиная с которого все приближения будут находиться в ε -окрестности решения.

В действительности же существует некоторая $\bar{\varepsilon}$ -окрестность решения, после попадания в которую приближения $\tilde{x}^{(n)}$ дальнейшего уточнения при выполнении итераций не происходит. В подтверждение сказанного приведем следующий результат.

Теорема 6.3. Пусть $\|B\| < 1$. Предположим, что вычисляемая на компьютере с помощью метода простых итераций последовательность $\tilde{x}^{(n)}$ удовлетворяет равенствам

$$\tilde{x}^{(n+1)} = B\tilde{x}^{(n)} + c + \xi^{(n)}, \quad n \geq 0, \quad (6.19)$$

где $\xi^{(n)} \in \mathbb{R}^m$, $\|\xi^{(n)}\| < \delta$. Если вычисления по формулам (6.6) и (6.19) начинаются с одного начального приближения $x^{(0)} = \tilde{x}^{(0)}$, то для всех $n \geq 1$ справедливы следующие оценки погрешности:

$$\|\tilde{x}^{(n)} - x^{(n)}\| \leq \bar{\varepsilon}, \quad (6.20)$$

$$\|\tilde{x}^{(n)} - \bar{x}\| \leq \|B\|^n \|\tilde{x}^{(0)} - \bar{x}\| + \bar{\varepsilon}, \quad (6.21)$$

$$\|\tilde{x}^{(n)} - \bar{x}\| \leq \frac{\|B\|}{1 - \|B\|} \|\tilde{x}^{(n)} - \tilde{x}^{(n-1)}\| + \bar{\varepsilon}; \quad (6.22)$$

здесь $\bar{\varepsilon} = \frac{\delta}{1 - \|B\|}$. ■

Таким образом, метод простых итераций устойчив, но гарантированная точность метода ограничена снизу величиной $\bar{\varepsilon}$. Критерий окончания (6.13) применим, если $\bar{\varepsilon} \ll \varepsilon$.

§ 6.2. Метод Зейделя

1. Описание метода. Пусть система (6.1) приведена к виду (6.4) с коэффициентами, вычисленными по формулам (6.5).

Метод Зейделя¹ можно рассматривать как модификацию метода Якоби. Основная идея модификации состоит в том, что при вычислении очередного $(k + 1)$ -го приближения к неизвестному x_i , при $i > 1$ используют уже найденные $(k + 1)$ -е приближения к неизвестным x_1, \dots, x_{i-1} , а не k -е приближения, как в методе Якоби.

На $(k + 1)$ -й итерации компоненты приближения $\mathbf{x}^{(k+1)}$ вычисляются по формулам

$$\begin{aligned} x_1^{(k+1)} &= b_{12}x_2^{(k)} + b_{13}x_3^{(k)} + \dots + b_{1m}x_m^{(k)} + c_1, \\ x_2^{(k+1)} &= b_{12}x_1^{(k+1)} + b_{23}x_3^{(k)} + \dots + b_{2m}x_m^{(k)} + c_2, \\ x_3^{(k+1)} &= b_{31}x_1^{(k+1)} + b_{32}x_2^{(k+1)} + \dots + b_{3m}x_m^{(k)} + c_3, \\ &\dots \\ x_m^{(k+1)} &= b_{m1}x_1^{(k+1)} + b_{m2}x_2^{(k+1)} + b_{m3}x_3^{(k+1)} + \dots + c_m. \end{aligned} \quad (6.23)$$

Введем нижнюю и верхнюю треугольные матрицы

$$\mathbf{B}_1 = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ b_{21} & 0 & 0 & \dots & 0 \\ b_{31} & b_{32} & 0 & \dots & 0 \\ \dots & & & & \\ b_{m1} & b_{m2} & b_{m3} & \dots & 0 \end{bmatrix}, \quad \mathbf{B}_2 = \begin{bmatrix} 0 & b_{12} & b_{13} & \dots & b_{1m} \\ 0 & 0 & b_{23} & \dots & b_{2m} \\ 0 & 0 & 0 & \dots & b_{3m} \\ \dots & & & & \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}.$$

Тогда расчетные формулы метода примут компактный вид:

$$\mathbf{x}^{(k+1)} = \mathbf{B}_1 \mathbf{x}^{(k+1)} + \mathbf{B}_2 \mathbf{x}^{(k)} + \mathbf{c}. \quad (6.24)$$

¹ Людвиг Зейдель (1821—1896) — немецкий астроном и математик.

Заметим, что $\mathbf{B} = \mathbf{B}_1 + \mathbf{B}_2$ и поэтому решение $\bar{\mathbf{x}}$ исходной системы удовлетворяет равенству

$$\bar{\mathbf{x}} = \mathbf{B}_1 \bar{\mathbf{x}} + \mathbf{B}_2 \bar{\mathbf{x}} + \mathbf{c}. \quad (6.25)$$

Метод Зейделя иногда называют также *методом Гаусса—Зейделя, процессом Либмана, методом последовательных замещений*.

2. Достаточные условия сходимости.

Теорема 6.4. Пусть $\|\mathbf{B}\| < 1$, где $\|\mathbf{B}\|$ — одна из норм $\|\mathbf{B}\|_\infty, \|\mathbf{B}\|_1$. Тогда при любом выборе начального приближения $\mathbf{x}^{(0)}$ метод Зейделя сходится со скоростью геометрической прогрессии, знаменатель которой $q \leq \|\mathbf{B}\|$. ■

Доказательство этой теоремы опускаем. Оно довольно громоздкое, хотя и не сложное. Приведем более компактное доказательство следующей теоремы, близкое к доказательству теоремы 6.1.

Теорема 6.5. Пусть выполнено условие

$$\|\mathbf{B}_1\| + \|\mathbf{B}_2\| < 1. \quad (6.26)$$

Тогда при любом выборе начального приближения метод Зейделя сходится и верна оценка погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| \leq q^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|, \quad (6.27)$$

где $q = \|\mathbf{B}_2\| / (1 - \|\mathbf{B}_1\|) < 1$.

Доказательство. Вычитая из равенства (6.24) равенство (6.25), имеем:

$$\mathbf{x}^{(k+1)} - \bar{\mathbf{x}} = \mathbf{B}_1(\mathbf{x}^{(k+1)} - \bar{\mathbf{x}}) + \mathbf{B}_2(\mathbf{x}^{(k)} - \bar{\mathbf{x}}). \quad (6.28)$$

Вычисляя нормы левой и правой частей этого равенства и используя свойства нормы, получаем:

$$\begin{aligned} \|\mathbf{x}^{(k+1)} - \bar{\mathbf{x}}\| &= \|\mathbf{B}_1(\mathbf{x}^{(k+1)} - \bar{\mathbf{x}}) + \mathbf{B}_2(\mathbf{x}^{(k)} - \bar{\mathbf{x}})\| \leq \\ &\leq \|\mathbf{B}_1\| \|\mathbf{x}^{(k+1)} - \bar{\mathbf{x}}\| + \|\mathbf{B}_2\| \|\mathbf{x}^{(k)} - \bar{\mathbf{x}}\|. \end{aligned}$$

Следовательно,

$$\|\mathbf{x}^{(k+1)} - \bar{\mathbf{x}}\| \leq q \|\mathbf{x}^{(k)} - \bar{\mathbf{x}}\|,$$

где $q = \|\mathbf{B}_2\| / (1 - \|\mathbf{B}_1\|)$.

Так как это неравенство верно для всех $k \geq 0$, то из него следует оценка (6.27). В силу условия (6.26) имеем $0 \leq q < 1$. Поэтому $\mathbf{x}^{(n)} \rightarrow \bar{\mathbf{x}}$ при $n \rightarrow \infty$. ■

Особо выделим часто встречающийся на практике случай систем с симметричными положительно определенными матрицами.

Теорема 6.6. Пусть A — симметричная положительно определенная матрица. Тогда при любом выборе начального приближения $x^{(0)}$ метод Зейделя сходится со скоростью геометрической прогрессии. ■

Отметим, что никаких дополнительных априорных условий типа малости нормы некоторой матрицы здесь не накладывается.

3. Апостериорная оценка погрешности

Предложение 6.2. Если выполнено условие $\|B\| < 1$, то для метода Зейделя справедлива апостериорная оценка погрешности

$$\|x^{(n)} - \bar{x}\| \leq \frac{\|B_2\|}{1 - \|B\|} \|x^{(n)} - x^{(n-1)}\|, \quad n \geq 1. \quad (6.29)$$

Доказательство. Положим $k = n - 1$ и запишем равенство (6.28) в следующем виде:

$$x^{(n)} - \bar{x} = B(x^{(n)} - \bar{x}) + B_2(x^{(n-1)} - x^{(n)}).$$

Тогда

$$\|x^{(n)} - \bar{x}\| \leq \|B\| \|x^{(n)} - \bar{x}\| + \|B_2\| \|x^{(n-1)} - x^{(n)}\|,$$

откуда и следует неравенство (6.29). ■

Полученное неравенство позволяет сформулировать простой критерий окончания итерационного процесса. Если требуется найти решение с точностью $\varepsilon > 0$, то итерации метода Зейделя следует вести

до выполнения неравенства $\frac{\|B_2\|}{1 - \|B\|} \|x^{(n-1)} - x^{(n)}\| < \varepsilon$ или эквивалентного ему неравенства

$$\|x^{(n)} - x^{(n-1)}\| < \varepsilon_2, \quad (6.30)$$

где $\varepsilon_2 = \frac{1 - \|B\|}{\|B_2\|} \varepsilon$.

4. Геометрическая интерпретация метода. Приведем геометрическую интерпретацию метода Зейделя при $m = 2$, т.е. когда он применяется для решения системы

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

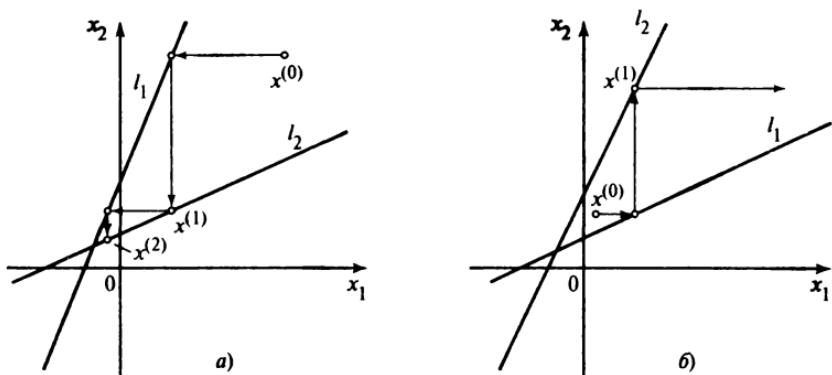


Рис. 6.1

Первое уравнение задает на плоскости x_1Ox_2 прямую l_1 , второе — прямую l_2 (рис. 6.1). Расчетные формулы метода принимают вид:

$$x_1^{(k+1)} = b_{12}x_2^{(k)} + c_1,$$

$$x_2^{(k+1)} = b_{21}x_1^{(k+1)} + c_2,$$

где $b_{12} = -a_{12}/a_{11}$, $c_1 = b_1/a_{11}$, $b_{21} = -a_{21}/a_{22}$, $c_2 = b_2/a_{22}$.

Пусть приближение $x^{(k)}$ уже найдено. Тогда при определении $x_1^{(k+1)}$ координата $x_2 = x_2^{(k)}$ фиксируется и точка x перемещается параллельно оси Ox_1 до пересечения с прямой l_1 . Координата x_1 точки пересечения принимается за $x_1^{(k+1)}$. Затем точка x перемещается вдоль прямой $x_1 = x_1^{(k+1)}$ до пересечения с прямой l_2 . Координата x_2 точки пересечения принимается за $x_2^{(k+1)}$.

На рис. 6.1, а, б приведены геометрические иллюстрации, отвечающие сходящемуся и расходящемуся итерационному процессу Зейделя. Видно, что характер сходимости может изменяться при перестановке уравнений.

Пример 6.2. Используя метод Зейделя, найдем решение системы (6.16) с точностью $\epsilon = 10^{-3}$.

После приведения системы к виду (6.17) убеждаемся, что $\|\mathbf{B}\|_\infty = 0.7278 < 1$ и поэтому в силу теоремы 6.4 метод Зейделя сходится.

Положим $x^{(0)} = (0, 0, 0)^T$ и будем вычислять последовательные приближения по формулам

$$\begin{aligned}x_1^{(k+1)} &= & 0.16x_2^{(k)} - 0.08x_3^{(k)} + 1.2, \\x_2^{(k+1)} &= & 0.2x_1^{(k+1)} - 0.424x_3^{(k)} - 1.736, \\x_3^{(k+1)} &= & -0.1389x_1^{(k+1)} - 0.5889x_2^{(k+1)} - 0.0667.\end{aligned}$$

Здесь

$$\mathbf{B}_2 = \begin{bmatrix} 0 & 0.16 & -0.08 \\ 0 & 0 & -0.424 \\ 0 & 0 & 0 \end{bmatrix}$$

и $\|\mathbf{B}_2\|_\infty = 0.424$. Будем вести итерации до выполнения критерия окончания (6.30), где $\varepsilon_2 = 10^{-3} \cdot (1 - 0.7278)/0.424 \approx 0.64 \cdot 10^{-3}$. Значения приближений с четырьмя цифрами после десятичной точки приведены в табл. 6.2.

Таблица 6.2

n	0	1	2	3	4
$x_1^{(n)}$	0.0000	1.2000	0.9088	0.8367	0.8121
$x_2^{(n)}$	0.0000	-1.4960	-1.8288	-1.9435	-1.9813
$x_3^{(n)}$	0.0000	0.6476	0.8841	0.9616	0.9873
$\ x^{(n)} - x^{(n-1)}\ _\infty$	—	1.4960	0.3328	0.1147	0.0378

Продолжение табл. 6.2

n	5	6	7	8
$x_1^{(n)}$	0.8040	0.8013	0.8004	0.8001
$x_2^{(n)}$	-1.9938	-1.9980	-1.9993	-1.9998
$x_3^{(n)}$	0.9958	0.9986	0.9995	0.9998
$\ x^{(n)} - x^{(n-1)}\ _\infty$	0.0125	0.0041	0.0014	0.0005

При $n = 8$ критерий окончания выполняется и можно положить $x_1 = -0.800 \pm 0.001$, $x_2 = -2.000 \pm 0.001$, $x_3 = 1.000 \pm 0.001$. Заметим, что в действительности решение с точностью $\varepsilon = 10^{-3}$ было получено уже при $n = 7$. ▲

5. Необходимое и достаточное условие сходимости. Метод Зейделя (6.24) может быть записан в эквивалентном виде

$$\mathbf{x}^{(k+1)} = \tilde{\mathbf{B}}\mathbf{x}^{(k)} + \tilde{\mathbf{c}},$$

где $\tilde{\mathbf{B}} = (\mathbf{E} - \mathbf{B}_1)^{-1}\mathbf{B}_2$, $\tilde{\mathbf{c}} = (\mathbf{E} - \mathbf{B}_1)^{-1}\mathbf{c}$. В такой форме записи он представляет собой метод простой итерации. Поэтому в силу теоремы 6.2 необходимым и достаточным условием сходимости метода является выполнение неравенства $\rho(\tilde{\mathbf{B}}) < 1$. Поскольку собственные числа матрицы $\tilde{\mathbf{B}}$ совпадают с числами λ , являющимися решениями обобщенной задачи на собственные значения $\mathbf{B}_2\mathbf{x} = \lambda(\mathbf{E} - \mathbf{B}_1)\mathbf{x}$, справедлива следующая теорема.

Теорема 6.7. *Метод Зейделя сходится при любом начальном приближении $\mathbf{x}^{(0)}$ тогда и только тогда, когда все решения обобщенной задачи на собственные значения $\mathbf{B}_2\mathbf{x} = \lambda(\mathbf{E} - \mathbf{B}_1)\mathbf{x}$ по модулю меньше единицы.* ■

Замечание. Существует устойчивое заблуждение, связанное с представлением о том, что метод Зейделя сходится быстрее, чем метод Якоби. Это действительно так, если матрица \mathbf{A} симметрична и положительно определена (мы убедились в преимуществе метода Зейделя для системы уравнений с такой матрицей, решив примеры 6.1 и 6.2). Однако в общем случае возможны ситуации, когда метод Якоби сходится, а метод Зейделя сходится медленнее или вообще расходится. Возможны и противоположные ситуации. Дело в том, что эти методы ориентированы на решение разных классов систем: метод Якоби — на системы с матрицами, близкими к диагональным, а метод Зейделя — на системы с матрицами, близкими к нижним треугольным.

§ 6.3. Метод последовательной верхней релаксации

Метод последовательной верхней релаксации является одним из наиболее эффективных и широко используемых итерационных методов для решения систем линейных алгебраических уравнений с симметричными положительно определенными матрицами \mathbf{A} . Этот метод часто называют SOR-методом¹. Частично популярность SOR-метода можно объяснить его простотой и тем, что он хорошо известен широкому кругу специалистов, занимающихся решением прикладных задач.

¹ От англ. successive over relaxation.

Суть *метода релаксации* состоит в следующем. После вычисления очередной i -й компоненты $(k+1)$ -го приближения по формуле метода Зейделя

$$\tilde{x}_i^{(k+1)} = b_{i1}x_1^{(k+1)} + b_{i2}x_2^{(k+1)} + \dots + b_{i,i-1}x_{i-1}^{(k+1)} + b_{i,i+1}x_{i+1}^{(k)} + \dots \\ \dots + b_{im}x_{im}^{(k)} + c_i$$

производят дополнительное смещение этой компоненты на величину $(\omega - 1)(\tilde{x}_i^{(k+1)} - x_i^{(k)})$, где ω — *параметр релаксации*. Таким образом, i -я компонента $(k+1)$ -го приближения вычисляется по формуле

$$x_i^{(k+1)} = \tilde{x}_i^{(k+1)} + (\omega - 1)(\tilde{x}_i^{(k+1)} - x_i^{(k)}) = \omega \tilde{x}_i^{(k+1)} + (1 - \omega)x_i^{(k)}.$$

На рис. 6.2 показано несколько первых итераций метода при значении параметра релаксации $\omega = 1.25$.

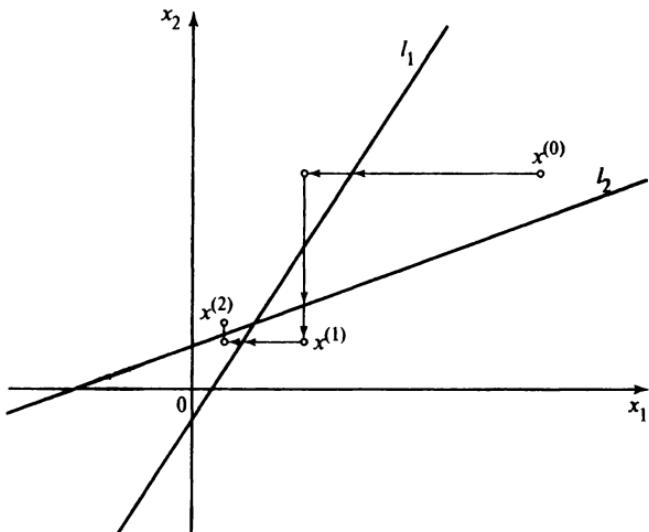


Рис. 6.2

В обозначениях, используемых в § 6.2, компактная формула для вычисления $x^{(k+1)}$ записывается следующим образом:

$$x^{(k+1)} = (1 - \omega)x^{(k)} + \omega B_1 x^{(k+1)} + \omega B_2 x^{(k)} + \omega c.$$

Как нетрудно видеть, при $\omega = 1$ метод релаксации совпадает с методом Зейделя. При $\omega > 1$ его было принято называть *методом последовательной верхней релаксации*, а при $\omega < 1$ — *методом последовательной нижней релаксации*. В последнее время метод релаксации

называют *методом последовательной верхней релаксации* для любых значений ω .

Если A — симметричная положительно определенная матрица, то при любом значении параметра ω ($0 < \omega < 2$) метод релаксации сходится. Часто оказывается возможным выбрать параметр $\omega > 1$ так, чтобы SOR-метод сходился существенно быстрее, чем методы Якоби и Зейделя. Однако выбор параметра релаксации — довольно трудная задача. Во многих случаях она решается экспериментальным путем.

Существуют различные модификации метода релаксации. Распространенный вариант метода связан с использованием различных параметров ω , для вычисления различных компонент x_i очередного $(k+1)$ -го приближения к решению.

Метод симметричной последовательной верхней релаксации (SSOR) заключается в чередовании итераций по методу SOR с такими же итерациями, осуществляемыми для системы, в которой уравнения занумерованы в обратном порядке.

Пример 6.3. Используя метод последовательной верхней релаксации с параметром $\omega = 1.12$, найдем решение системы (6.16) с точностью $\epsilon = 10^{-3}$.

Приведем систему к виду (6.17), положим $x^{(0)} = (0, 0, 0)^T$ и будем вычислять последовательные приближения по формулам:

$$x_1^{(k+1)} = (1 - \omega)x_1^{(k)} + \omega(0.16x_2^{(k)} - 0.08x_3^{(k)} + 1.2),$$

$$x_2^{(k+1)} = \omega \cdot 0.2x_1^{(k+1)} + (1 - \omega)x_2^{(k)} + \omega(-0.424x_3^{(k)} - 1.736),$$

$$x_3^{(k+1)} = \omega(-0.1389x_1^{(k+1)} - 0.5889x_2^{(k+1)}) + (1 - \omega)x_3^{(k)} - \omega \cdot 0.0667.$$

Значения приближений с четырьмя цифрами после десятичной точки приведены в табл. 6.3.

Таблица 6.3

n	0	1	2	3	4	5
$x_1^{(n)}$	0 0000	1 3440	0 8166	0 8094	0 7995	0 8001
$x_2^{(n)}$	0 0000	-1 6433	-1 9442	-1 9973	-1 9998	-2 0000
$x_3^{(n)}$	0 0000	0 8001	0 9846	0 9986	1 0001	1 0000

Сравнение с точным решением $x_1 = 0.8$, $x_2 = -2$, $x_3 = 1$ показывает, что для получения приближения к решению с точностью $\epsilon = 10^{-3}$ потребовалось всего четыре итерации. Напомним, что для достижения той же точности при том же начальном приближении методами Якоби и Зейделя потребовалось соответственно тринадцать и семь итераций.

§ 6.4. Другие двухслойные итерационные методы

Наряду с методами простой итерации, Зейделя и последовательной верхней релаксации существует большое число иных итерационных методов, предназначенных для решения систем линейных алгебраических уравнений

$$Ax = b. \quad (6.31)$$

Особенно глубоко разработана теория итерационных методов для систем с симметричными положительно определенными матрицами. Рассмотрим некоторые из них.

Всюду в этом параграфе будем считать, что матрица A — симметричная и положительно определенная. Напомним, что для такой матрицы все ее собственные значения $\lambda_1, \dots, \lambda_m$ положительны. Кроме того, справедливы неравенства

$$\lambda_{\min} \|x\|_2^2 \leq (Ax, x) \leq \lambda_{\max} \|x\|_2^2 \quad \text{для всех } x \in \mathbb{R}^m, \quad (6.32)$$

$$\lambda_{\min} \|x\|_2 \leq \|Ax\|_2 \leq \lambda_{\max} \|x\|_2 \quad \text{для всех } x \in \mathbb{R}^m, \quad (6.33)$$

где λ_{\min} и λ_{\max} — минимальное и максимальное собственные значения.

Из неравенств (6.32), в частности, следует, что величина $\|x\|_A = \sqrt{(Ax, x)}$ является нормой в \mathbb{R}^m ; эту норму принято называть *энергетической*. Неравенства (6.33) позволяют увидеть, что $\|A\|_2 = \lambda_{\max}$, $\|A^{-1}\|_2 = \lambda_{\min}^{-1}$ и поэтому $\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 = \lambda_{\max} / \lambda_{\min}$.

1. Двухслойные итерационные методы. Итерационный метод называют *одношаговым* или *двухслойным* если для вычисления очередного приближения $x^{(k+1)}$ к решению \bar{x} системы (6.31) используется только одно предыдущее приближение $x^{(k)}$. Например, методы простой итерации, Зейделя и последовательной верхней релаксации — двухслойные.

Многие двухслойные итерационные методы могут быть записаны в каноническом виде

$$B \frac{x^{(k+1)} - x^{(k)}}{\tau_{k+1}} + Ax^{(k)} = b, \quad k = 0, 1, \dots . \quad (6.34)$$

Здесь B — некоторая невырожденная матрица, $\tau_{k+1} > 0$ — *итерационные параметры*. Матрица B может, вообще говоря, зависеть от k , то есть меняться от итерации к итерации. В случае, когда параметры $\tau_{k+1} = \tau$ (и матрица B) не зависят от номера итерации, метод называют *стационарным*.

Если $B = E$ — единичная матрица, итерационный метод принимает вид

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

и называется явным, поскольку в нем очередное приближение явным образом выражается через предыдущее:

$$x^{(k+1)} = x^{(k)} - \tau_{k+1} (Ax^{(k)} - b), \quad k = 0, 1, \dots \quad (6.35)$$

В общем случае, при $B \neq E$, метод является неявным, так как на каждой итерации требуется решать систему уравнений

$$Bx^{(k+1)} = Bx^{(k)} - \tau_{k+1} (Ax^{(k)} - b), \quad k = 0, 1, \dots$$

Конечно, матрица B должна быть такой, чтобы каждая из этих систем решалась значительно быстрее, чем исходная система (6.31). Например, матрица B может быть диагональной, трехдиагональной, треугольной или разреженной.

Замечание. При численной реализации неявного метода (6.34) обычно сначала вычисляется невязка $r^{(k)} = b - Ax^{(k)}$, затем относительно вектора $w^{(k)} = \frac{x^{(k+1)} - x^{(k)}}{\tau_{k+1}}$, пропорционального поправке к решению, решается система уравнений

$$Bw^{(k)} = r^{(k)}, \quad (6.36)$$

после чего находится очередное приближение $x^{(k+1)} = x^{(k)} + \tau_{k+1} w^{(k)}$.

2. Каноническая форма записи методов Якоби, Зейделя и последовательной верхней релаксации. Представим матрицу A в виде суммы $A = D + A_1 + A_2$, где D — диагональная матрица с диагональными элементами $d_{ii} = a_{ii}$ для $1 \leq i \leq m$, матрица A_1 — нижняя треугольная с элементами $a_{ij}^{(1)} = a_{ij}$ для $1 \leq j < i \leq m$ и $a_{ij}^{(1)} = 0$ для $1 \leq i < j \leq m$; наконец, матрица A_2 — верхняя треугольная с элементами $a_{ij}^{(2)} = a_{ij}$ для $1 \leq i < j \leq m$ и $a_{ij}^{(2)} = 0$ для $1 \leq j < i \leq m$.

Рассмотренные в предыдущих параграфах методы могут быть записаны в каноническом виде (6.34). Метод Якоби соответствует выбору $B = D$ и $\tau_{k+1} = 1$. Метод Зейделя может быть записан в виде (6.34), если взять $B = D + A_1$ (т.е. выбрать в качестве B нижнюю треугольную часть

матрицы A) и $\tau_{k+1} = 1$. Метод последовательной верхней релаксации соответствует выбору $B = D + \omega A_1$ и $\tau_{k+1} = \omega$.

3. Явный стационарный итерационный метод с оптимальным параметром. Рассмотрим явный стационарный итерационный метод

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau(A\mathbf{x}^{(k)} - \mathbf{b}), \quad k = 0, 1, \dots. \quad (6.37)$$

Как нетрудно видеть, он представляет собой метод простой итерации

$$\mathbf{x}^{(k+1)} = S_\tau \mathbf{x}^{(k)} + \tau \mathbf{b}, \quad k = 0, 1, \dots,$$

где $S_\tau = E - \tau A$ — матрица, называемая *матрицей перехода*. Эта матрица симметрична, причем ее собственные значения $\lambda_i(S_\tau)$ вычисляются через собственные значения матрицы A по формуле $\lambda_i(S_\tau) = 1 - \tau \lambda_i$. Поэтому спектральный радиус матрицы S_τ совпадает с $\|S_\tau\|_2 = \|E - \tau A\|_2$ и вычисляется по формуле $\rho(S_\tau) = q(\tau) \equiv \max_{1 \leq i \leq m} |1 - \tau \lambda_i|$.

В силу теоремы 6.2 метод (6.37) сходится тогда и только тогда, когда $q(\tau) < 1$, причем справедлива оценка погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq q(\tau)^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2, \quad n \geq 1.$$

Скорость сходимости метода зависит от τ и она тем выше, чем меньше величина $q(\tau)$. Можно показать, что функция $q(\tau)$ принимает свое минимальное значение q_0 при

$$\tau = \tau_{\text{опт}} = \frac{2}{\lambda_{\min} + \lambda_{\max}}. \quad (6.38)$$

(В этом смысле выбор параметра (6.38) является оптимальным.) При этом

$$q_0 = q(\tau_{\text{опт}}) = \min_{0 < \tau} q(\tau) = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} = \frac{\lambda_{\max}/\lambda_{\min} - 1}{\lambda_{\max}/\lambda_{\min} + 1},$$

т.е.

$$q_0 = \frac{\text{cond}_2(A) - 1}{\text{cond}_2(A) + 1}. \quad (6.39)$$

Таким образом, справедлив следующий результат.

Теорема 6.8. Явный стационарный итерационный метод (6.37) с оптимальным выбором параметра (6.38) сходится с оценкой погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq q_0^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2, \quad n \geq 1, \quad (6.40)$$

где q_0 определяется формулой (6.39). ■

Замечание 1. Оценка (6.40) позволяет оценить число итераций, достаточное для того, чтобы уменьшить начальную погрешность в нужное число раз, т.е. добиться выполнения неравенства

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq \varepsilon \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2 \quad (6.41)$$

при фиксированном $\varepsilon > 0$. Это неравенство выполнено, если $q_0^n \leq \varepsilon$, т.е. если

$$n \geq n_0(\varepsilon) = \frac{\ln(1/\varepsilon)}{\ln(1/q_0)}.$$

Замечание 2. Из формулы (6.39) видно, что q_0 растет с ростом $\text{cond}_2(\mathbf{A})$, причем $q_0 = 1 - \frac{2}{\text{cond}_2(\mathbf{A}) + 1} \rightarrow 1$ при $\text{cond}_2(\mathbf{A}) \rightarrow \infty$. Таким образом, скорость сходимости метода с ростом числа обусловленности матрицы \mathbf{A} замедляется. Кроме того,

$$n_0(\varepsilon) \approx \frac{1}{2} \text{cond}_2(\mathbf{A}) \ln(1/\varepsilon) \quad \text{при } \text{cond}_2(\mathbf{A}) \gg 1. \quad (6.42)$$

Следует обратить внимание на то, что в прикладных задачах часто встречаются системы уравнений с матрицами, для которых число обусловленности $\text{cond}_2(\mathbf{A}) \gg 1$. В таких задачах число итераций, определяемое по формуле (6.42), недопустимо велико.

Замечание 3. Выбор оптимального параметра $\tau_{\text{опт}}$ предполагает знание собственных значений λ_{\min} и λ_{\max} . Часто они неизвестны, но имеются их оценки $0 < \mu \leq \lambda_{\min} \leq \lambda_{\max} \leq M$. В этой ситуации

оптимальным является выбор параметра $\tau_{\text{опт}} = \frac{2}{\mu + M}$, а метод сходится с оценкой (6.38), в которой $q_0 = \frac{M - \mu}{M + \mu} = \frac{M/\mu - 1}{M/\mu + 1}$. Кроме того,

$$n_0(\varepsilon) \approx \frac{1}{2} (M/\mu) \ln(1/\varepsilon) \quad \text{при } M/\mu \gg 1. \quad (6.43)$$

4. Явный итерационный метод с чебышевским набором параметров. Заметим, что погрешность явного итерационного метода

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau_{k+1}(\mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}), \quad k = 0, 1, \dots \quad (6.44)$$

удовлетворяет равенству

$$\mathbf{x}^{(k+1)} - \bar{\mathbf{x}} = S_{\tau_{k+1}}(\mathbf{x}^{(k)} - \bar{\mathbf{x}}), \quad k \geq 0,$$

где $S_{\tau_{k+1}} = E - \tau_{k+1}A$. Поэтому для погрешности, полученной на n -й итерации справедливо равенство

$$\mathbf{x}^{(n)} - \bar{\mathbf{x}} = S_{\tau_n} S_{\tau_{n-1}} \dots S_{\tau_1} (\mathbf{x}^{(0)} - \bar{\mathbf{x}})$$

и, как следствие, — оценка

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq q_n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2,$$

где $q_n = \|S_{\tau_n} S_{\tau_{n-1}} \dots S_{\tau_1}\|_2$.

Зададимся фиксированным числом итераций n и поставим задачу выбора таких итерационных параметров $\tau_1, \tau_2, \dots, \tau_n$, для которых величина q_n минимальна. Известно, что такими параметрами являются

$$\tau_k = \frac{\tau_{\text{опт}}}{1 + q_0 \mu_k}, \quad k = 1, 2, \dots, n, \quad (6.45)$$

где $\tau_{\text{опт}}$ и q_0 задаются формулами (6.38) и (6.39), а числа

$$\mu_k = \cos \frac{j_k \pi}{2n}, \quad k = 1, 2, \dots, n \quad (6.46)$$

— это занумерованные в некотором порядке корни многочлена Чебышева T_n^1 . Здесь $J_n = (j_1, j_2, \dots, j_n)$ — упорядоченный набор нечетных чисел из множества $\{1, 3, \dots, 2n - 1\}$, определяющий порядок, в котором берутся корни.

Метод (6.44) с итерационными параметрами (6.45) называется явным итерационным методом с чебышевским набором параметров или методом Ричардсона.

Теорема 6.9. Для явного итерационного метода с чебышевским набором параметров справедлива оценка погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq \frac{2\rho^n}{1 + \rho^{2n}} \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2, \quad (6.47)$$

в которой $\rho = \frac{\sqrt{\text{cond}_2(A)} - 1}{\sqrt{\text{cond}_2(A)} + 1}$. ■

Используя оценку (6.47), определим число итераций, достаточное для того, чтобы добиться выполнения неравенства (6.41), т.е. уменьшить начальную погрешность в $1/\epsilon$ раз. Заметим, что неравенство

¹ Определение и свойства многочленов Чебышева рассматриваются в § 11.6.

$\frac{2\rho^n}{1+\rho^{2n}} \leq \varepsilon$ выполнено, если $\rho^n \leq \varepsilon/2$. Отсюда $n \geq n_0(\varepsilon) = \frac{\ln(2/\varepsilon)}{\ln(1/\rho)}$.

Поскольку $\rho = 1 - \frac{2}{\sqrt{\text{cond}_2(A)} + 1}$, то

$$n_0(\varepsilon) \approx \frac{1}{2} \sqrt{\text{cond}_2(A)} \ln(2/\varepsilon) \quad \text{при } \text{cond}_2(A) \gg 1. \quad (6.48)$$

Сравнение значения (6.48) со значением $n_0(\varepsilon)$, определяемым формулой (6.42) показывает, что метод с чебышевским набором параметров требует примерно в $\sqrt{\text{cond}_2(A)}$ меньшее число итераций, т.е. он обладает существенным преимуществом в числе итераций перед методом простой итерации с оптимальным выбором параметра.

Замечание 1. Следует обратить внимание на то, что число итераций n для итерационного метода с чебышевским набором параметров заранее фиксируется и за приближение к решению \bar{x} принимается именно $x^{(n)}$. Значения $x^{(k)}$, полученные на промежуточных шагах с $0 < k < n$, плохо аппроксимируют решение. При необходимости уточнить приближение цикл из n итераций повторяется, и за следующее приближение к решению принимается $x^{(2n)}$.

Замечание 2. При $n = 1$ метод с чебышевским набором параметров превращается в явный стационарный итерационный метод (6.37) с оптимальным параметром (6.38).

Замечание 3. Метод Ричардсона длительное время не использовался на практике из-за серьезных проблем с вычислительной устойчивостью. Как оказалось, для вычислительной устойчивости рассматриваемого метода очень важен порядок, в котором берутся нули многочлена Чебышева. Устойчивость метода имеет место только при специальной нумерации нулей. В противном случае наблюдается резкий рост вычислительной погрешности, который может привести к полной потере точности.

В настоящее время известны наборы чисел¹ $J_n = (j_1, j_2, \dots, j_n)$, обеспечивающие вычислительную устойчивость метода, и алгоритмы их построения для любого n . Приведем такие наборы для $n = 4, 8, 16$.

$$J_4 = (1, 7, 3, 5), \quad J_8 = (1, 15, 7, 9, 3, 13, 5, 11),$$

$$J_{16} = (1, 31, 15, 17, 7, 25, 9, 23, 3, 29, 13, 19, 5, 27, 11, 21).$$

¹ Алгоритмы их построения были опубликованы практически одновременно в статье В.А. Лебедев, С.А. Финогенов // Журн. вычисл. матем. и матем. физ. 1971. Т. 11. № 2. С. 425—438 и в монографии А.А. Самарский. Теория разностных схем. М.: Наука, 1971. Дополнительную информацию можно найти также в [62, 89].

Замечание 4. Если вместо собственных значений λ_{\min} и λ_{\max} известны лишь их оценки $0 < \mu \leq \lambda_{\min} \leq \lambda_{\max} \leq M$, то чебышевский набор параметров может быть вычислен по формулам (6.45), (6.46) с $\tau_{\text{опт}} = \frac{2}{\mu + M}$, $q_0 = \frac{M - \mu}{M + \mu}$. Для построенного таким образом метода будет справедлива оценка (6.47), в которой $\rho = \frac{\sqrt{M/\mu} - 1}{\sqrt{M/\mu} + 1}$. Кроме того, число итераций, достаточное для уменьшения погрешности в $1/\varepsilon$ раз составит

$$n_0(\varepsilon) \approx \frac{1}{2} \sqrt{M/\mu} \ln(1/\varepsilon) \quad \text{при } M/\mu \gg 1$$

(ср. с (6.43)).

6. Неявный итерационный метод с оптимальным выбором параметра. Пусть B — симметричная положительно определенная матрица. Известно, что существует единственная симметричная положительно определенная матрица C такая¹, что $C^2 = B$. Положим $y^{(k)} = Cx^{(k)}$ и запишем неявный итерационный метод (6.34) в виде эквивалентного явного метода

$$\frac{y^{(k+1)} - y^{(k)}}{\tau_{k+1}} + \tilde{A}y^{(k)} = \tilde{b}, \quad k = 0, 1, \dots,$$

где $\tilde{A} = C^{-1}AC^{-1}$, $\tilde{b} = C^{-1}b$.

Поскольку матрицы A и C^{-1} — симметричные и положительно определенные, то и матрица \tilde{A} — симметричная и положительно определенная. Обозначим через $\tilde{\lambda}_{\min}$ и $\tilde{\lambda}_{\max}$ ее минимальное и максимальное собственные значения. Заметим, что $\bar{y} = Cx$ является решением системы $\tilde{A}y = \tilde{b}$.

Из теоремы 6.8 следует, что явный стационарный итерационный метод

$$\frac{y^{(k+1)} - y^{(k)}}{\tau} + \tilde{A}y^{(k)} = \tilde{b}, \quad k = 0, 1, \dots,$$

с оптимальным выбором параметра

$$\tau = \tau_{\text{опт}} = \frac{2}{\tilde{\lambda}_{\min} + \tilde{\lambda}_{\max}} \tag{6.49}$$

¹ Матрица C называется *корнем квадратным из матрицы A* и обозначается через \sqrt{A} .

сходится с оценкой погрешности

$$\|\mathbf{y}^{(n)} - \bar{\mathbf{y}}\|_2 \leq \tilde{q}_0^n \|\mathbf{y}^{(0)} - \bar{\mathbf{y}}\|_2, \quad n \geq 1,$$

в которой

$$\tilde{q}_0 = \frac{\operatorname{cond}_2(\tilde{\mathbf{A}}) - 1}{\operatorname{cond}_2(\tilde{\mathbf{A}}) + 1}. \quad (6.50)$$

Как следствие, эквивалентный неявный стационарный метод

$$\mathbf{B} \frac{\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}}{\tau} + \mathbf{A}\mathbf{x}^{(k)} = \mathbf{b}, \quad k = 0, 1, \dots \quad (6.51)$$

с оптимальным выбором параметра (6.49) сходится с оценкой погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq \sqrt{\operatorname{cond}_2(\mathbf{B})} \tilde{q}_0^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2, \quad n \geq 1, \quad (6.52)$$

в которой \tilde{q}_0 определяется формулой (6.50).

Для того, чтобы этот метод обладал преимуществом в скорости сходимости перед явным методом (6.37) с оптимальным выбором параметра, нужно, чтобы число \tilde{q}_0 было существенно меньше числа q_0 , т.е. нужно выполнение условия $\operatorname{cond}_2(\tilde{\mathbf{A}}) \ll \operatorname{cond}_2(\mathbf{A})$.

Матрицу \mathbf{B} , гарантирующую существенное уменьшение числа обусловленности, часто называют *предобусловливателем*¹, а процесс ее построения — *предобусловливанием*. Предобусловливание является весьма сложным делом, поскольку предобусловливатель \mathbf{B} обязан обладать еще одним важным свойством: он должен иметь такую структуру, чтобы системы вида (6.36) решались достаточно быстро.

Как правило, найти собственные значения матрицы $\tilde{\mathbf{A}} = \mathbf{C}^{-1}\mathbf{A}\mathbf{C}^{-1}$ являющиеся решениями обобщенной задачи на собственные значения $\mathbf{A}\mathbf{x} = \lambda\mathbf{B}\mathbf{x}$, весьма сложно. Значительно проще найти оценки μ_1 и M_1 минимального и максимального собственных значений: $0 < \mu_1 \leq \tilde{\lambda}_{\min} \leq \tilde{\lambda}_{\max} \leq M_1$ или, что то же самое, — найти числа $0 < \mu_1 \leq M_1$ такие, что²

$$\mu_1(\mathbf{B}\mathbf{x}, \mathbf{x}) \leq (\mathbf{A}\mathbf{x}, \mathbf{x}) \leq M_1(\mathbf{B}\mathbf{x}, \mathbf{x}) \quad \text{для всех } \mathbf{x} \in \mathbb{R}^m. \quad (6.53)$$

¹ Иногда ее называют *переобусловливателем*.

² Матрицы \mathbf{B} и \mathbf{A} , связанные неравенством (6.53), называют *спектрально эквивалентными* или *энергетически эквивалентными*.

В этом случае можно использовать итерационный метод (6.51) с

$$\tau = \frac{2}{\mu_1 + M_1}, \quad (6.54)$$

для которого справедлива оценка погрешности (6.52) с

$$\tilde{q}_0 = \frac{M_1 - \mu_1}{M_1 + \mu_1} = \frac{M_1/\mu_1 - 1}{M_1/\mu_1 + 1}. \quad (6.55)$$

Замечание. Следует подчеркнуть, что для реализации метода (6.51) нет необходимости в вычислении матриц C и \tilde{A} . Они нужны лишь для теоретического исследования свойств метода.

6. Неявный итерационный метод с чебышевским набором параметров. Аналогичные соображения приводят к пониманию того, что в неявном методе (6.34) с симметричной положительно определенной матрицей B возможен выбор чебышевского набора параметров

$$\tau_k = \frac{\tau_{\text{опт}}}{1 + \tilde{q}_0 \mu_k}, \quad k = 1, 2, \dots, n, \quad (6.56)$$

в котором $\tau_{\text{опт}}$, \tilde{q}_0 и μ_k вычисляются по формулам (6.49), (6.50) и (6.46). Этот метод представляет собой *неявный итерационный метод с чебышевским набором параметров*. Для него справедлива оценка погрешности

$$\|x^{(n)} - \bar{x}\|_2 \leq \sqrt{\text{cond}_2(B)} \frac{2\tilde{\rho}^n}{1 + \tilde{\rho}^{2n}} \|x^{(0)} - \bar{x}\|_2, \quad (6.57)$$

$$\text{в которой } \tilde{\rho} = \frac{\sqrt{\text{cond}_2(\tilde{A})} - 1}{\sqrt{\text{cond}_2(\tilde{A})} + 1}.$$

Если вместо значений $\tilde{\lambda}_{\min}$ и $\tilde{\lambda}_{\max}$ известны лишь их оценки $0 < \mu_1 \leq \tilde{\lambda}_{\min} \leq \tilde{\lambda}_{\max} \leq M_1$, удовлетворяющие неравенству (6.53), то в формуле

(6.56) следует взять $\tau_{\text{опт}} = \frac{2}{\mu_1 + M_1}$, $\tilde{q}_0 = \frac{M_1 - \mu_1}{M_1 + \mu_1}$. Для метода (6.34) с таким набором параметров будет справедлива оценка погрешности

$$(6.57) \text{ с } \tilde{\rho} = \frac{\sqrt{M_1/\mu_1} - 1}{\sqrt{M_1/\mu_1} + 1}.$$

Конечно, при реализации неявного метода с чебышевским набором параметров необходимо использовать порядок выбора нулей многочлена Чебышева, гарантирующий вычислительную устойчивость.

7. Методы минимальных невязок и поправок. Запишем двухслойный итерационный метод (6.35) в виде

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \tau_{k+1} \mathbf{r}^{(k)}, \quad k \geq 0, \quad (6.58)$$

где $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$ — невязка, полученная на k -й итерации. Умножив слева обе части равенства (6.58) на $-\mathbf{A}$, получим соотношение

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \tau_{k+1} \mathbf{A}\mathbf{r}^{(k)}. \quad (6.59)$$

Выберем параметр τ_{k+1} таким образом, чтобы минимизировать евклидову норму невязки $\|\mathbf{r}^{(k+1)}\|_2 = \|\mathbf{r}^{(k)} - \tau_{k+1} \mathbf{A}\mathbf{r}^{(k)}\|_2$. Заметим, что функция

$\varphi(\tau_{k+1}) \equiv \|\mathbf{r}^{(k)} - \tau_{k+1} \mathbf{A}\mathbf{r}^{(k)}\|_2^2 = \|\mathbf{r}^{(k)}\|_2^2 - 2\tau_{k+1} (\mathbf{A}\mathbf{r}^{(k)}, \mathbf{r}^{(k)}) + \tau_{k+1}^2 \|\mathbf{A}\mathbf{r}^{(k)}\|_2^2$ является квадратичной по τ_{k+1} и достигает своего минимума при значении

$$\tau_{k+1} = \frac{(\mathbf{A}\mathbf{r}^{(k)}, \mathbf{r}^{(k)})}{\|\mathbf{A}\mathbf{r}^{(k)}\|_2^2}. \quad (6.60)$$

Теорема 6.10. Для метода минимальных невязок (6.58)–(6.60) справедлива оценка

$$\|\mathbf{r}^{(n)}\|_2 \leq q_0^n \|\mathbf{r}^{(0)}\|_2, \quad n \geq 1, \quad (6.61)$$

где q_0 определяется формулой (6.39). Как следствие, метод сходится с оценкой погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq \text{cond}_2(\mathbf{A}) q_0^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2, \quad n \geq 0. \quad (6.62)$$

Доказательство. В силу специального выбора параметра τ_{k+1} справедливо неравенство

$$\|\mathbf{r}^{(k+1)}\|_2 = \|(\mathbf{E} - \tau_{k+1} \mathbf{A})\mathbf{r}^{(k)}\|_2 = \min_{\tau > 0} \|(\mathbf{E} - \tau \mathbf{A})\mathbf{r}^{(k)}\|_2 \leq \|(\mathbf{E} - \tau_{\text{опт}} \mathbf{A})\mathbf{r}^{(k)}\|_2.$$

Принимая во внимание, что $\|\mathbf{E} - \tau_{\text{опт}} \mathbf{A}\|_2 = q_0$ (см. п. 3), имеем

$$\|\mathbf{r}^{(k+1)}\|_2 \leq \|\mathbf{E} - \tau_{\text{опт}} \mathbf{A}\|_2 \|\mathbf{r}^{(k)}\|_2 = q_0 \|\mathbf{r}^{(k)}\|_2, \quad k \geq 0.$$

Как следствие, справедлива оценка (6.61).

Поскольку $\|\mathbf{r}^{(n)}\|_2 = \|\mathbf{A}(\mathbf{x}^{(n)} - \bar{\mathbf{x}})\|_2$, то

$$\lambda_{\min} \|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq \|\mathbf{r}^{(n)}\|_2, \quad \|\mathbf{r}^{(0)}\|_2 \leq \lambda_{\max} \|\mathbf{x}^{(0)}\|_2$$

и неравенство (6.62) следует из неравенства (6.61). ■

Замечание Метод минимальных невязок сходится с той же скоростью, что и явный итерационный метод с оптимальным выбором параметра. Однако он обладает тем существенным преимуществом, что не требует для своей реализации знания собственных значений λ_{\min} и λ_{\max} или их оценок.

Неявный итерационный метод

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \tau_{k+1} \mathbf{w}^{(k)}, \quad \mathbf{B}\mathbf{w}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}, \quad k = 0, 1, \dots,$$

являющийся аналогом метода минимальных невязок, называется *методом минимальных поправок*. В этом методе итерационные параметры вычисляются по формуле

$$\tau_{k+1} = \frac{(\mathbf{A}\mathbf{w}^{(k)}, \mathbf{w}^{(k)})}{(\mathbf{B}^{-1}\mathbf{A}\mathbf{w}^{(k)}, \mathbf{A}\mathbf{w}^{(k)})}, \quad k = 0, 1, \dots,$$

а очередное значение поправки по формуле

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \tau_{k+1} \mathbf{B}^{-1} \mathbf{A} \mathbf{w}^{(k)}$$

Если матрицы \mathbf{A} и \mathbf{B} спектрально эквивалентны, т.е. связаны неравенством (6.53) с $0 < \mu_1 < M_1$, то для метода минимальных поправок справедлива оценка

$$\|\mathbf{r}^{(k+1)}\|_{\mathbf{B}^{-1}} \leq \tilde{q}_0 \|\mathbf{r}^{(k)}\|_{\mathbf{B}^{-1}}, \quad k = 0, 1, \dots,$$

где $\tilde{q}_0 = \frac{M_1 - \mu_1}{M_1 + \mu_1}$. Как следствие, справедлива оценка погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq c \tilde{q}_0^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2, \quad n \geq 1,$$

где $c = \sqrt{\text{cond}_2(\mathbf{B})} (M_1 / \mu_1)$

8. Метод наискорейшего спуска. Покажем, что задача решения системы уравнений (6.31) с симметричной положительно определенной матрицей \mathbf{A} эквивалентна задаче поиска точки минимума квадратичной функции

$$F(\mathbf{x}) = \frac{1}{2} (\mathbf{A}\mathbf{x}, \mathbf{x}) - (\mathbf{b}, \mathbf{x}). \quad (6.63)$$

Действительно, пусть $\bar{\mathbf{x}}$ — решение системы (6.31). Тогда для $\mathbf{x} = \bar{\mathbf{x}} + \Delta\mathbf{x}$ при $\Delta\mathbf{x} \neq 0$ имеем

$$\begin{aligned} F(\mathbf{x}) - F(\bar{\mathbf{x}}) &= \frac{1}{2} (\mathbf{A}(\bar{\mathbf{x}} + \Delta\mathbf{x}), \bar{\mathbf{x}} + \Delta\mathbf{x}) - (\mathbf{b}, \mathbf{x}) - \frac{1}{2} (\mathbf{A}\bar{\mathbf{x}}, \bar{\mathbf{x}}) + (\mathbf{b}, \bar{\mathbf{x}}) = \\ &= (\mathbf{A}\bar{\mathbf{x}} - \mathbf{b}, \Delta\mathbf{x}) + \frac{1}{2} (\mathbf{A}\Delta\bar{\mathbf{x}}, \Delta\bar{\mathbf{x}}) = \frac{1}{2} (\mathbf{A}\Delta\bar{\mathbf{x}}, \Delta\bar{\mathbf{x}}) > 0. \end{aligned}$$

Таким образом, всякий итерационный метод минимизации квадратичной функции (6.63) одновременно является итерационным методом решения системы (6.31). Например, метод покоординатного спуска, примененный к функции (6.63), дает ту же последовательность приближений, что и метод Зейделя. Подробнее об этом будет сказано в § 10.2; см. также [9].

В методе наискорейшего спуска итерации осуществляются по тем же формулам

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \tau_{k+1} \mathbf{r}^{(k)}, \quad \mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \tau_{k+1} \mathbf{A} \mathbf{r}^{(k)}, \quad k \geq 0, \quad (6.64)$$

что и в методе минимальных невязок. Но параметр τ_{k+1} вычисляется иначе:

$$\tau_{k+1} = \frac{\|\mathbf{r}^{(k)}\|_2^2}{(\mathbf{A}\mathbf{r}^{(k)}, \mathbf{r}^{(k)})}. \quad (6.65)$$

Именно при этом значении достигает минимума функция $\varphi_k(\alpha) = F(\mathbf{x}^{(k)} + \alpha \mathbf{r}^{(k)})$ (см. § 10.3).

Теорема 6.11. *Метод наискорейшего спуска (6.64), (6.65) сходится с оценкой погрешности*

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_A \leq q_0^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_A, \quad n \geq 1, \quad (6.66)$$

где q_0 определяется формулой (6.39). Как следствие,

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\|_2 \leq \sqrt{\text{cond}_2(\mathbf{A})} q_0^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|_2, \quad n \geq 1. \blacksquare$$

Замечание 1. Как и метод минимальных невязок, метод наискорейшего спуска сходится с той же скоростью, что и явный итерационный метод с оптимальным выбором параметра. Однако он обладает тем существенным преимуществом перед явным итерационным методом с оптимальным выбором параметра, что не требует для своей реализации знания собственных значений λ_{\min} и λ_{\max} или их оценок.

Замечание 2. Дополнительную информацию о методе наискорейшего спуска можно найти в § 10.3.

§ 6.5. Метод сопряженных градиентов

1. Метод сопряженных градиентов.¹ Совместная статья Хестенса и Штифеля² появилась в 1952 году после того, как они независимо открыли метод сопряженных градиентов.

Первоначально этот метод рассматривался как один из прямых методов решения систем линейных алгебраических уравнений

$$Ax = b$$

с самосопряженными положительно определенными матрицами порядка m .

Было доказано, что в точной арифметике этот метод дает точное решение за m шагов. Однако на практике точная сходимость за m шагов не наблюдалась, поэтому метод сопряженных градиентов расценили как неустойчивый и интерес к нему на многие годы был утрачен. Только после того, как в 1971 году Рейд³ продемонстрировал значимость метода сопряженных градиентов как итерационного метода решения систем с разреженными матрицами, началось его широкое применение. В настоящее время метод сопряженных градиентов по праву рассматривается как один из наиболее эффективных итерационных методов решения систем уравнений с самосопряженными положительно определенными матрицами.

Рассмотрим этот метод. Пусть $x^{(0)}$ — начальное приближение к решению, $r^{(0)} = b - Ax^{(0)}$ — начальное значение невязки и $p^{(0)} = r^{(0)}$. Приближения $x^{(k+1)}$ и невязки $r^{(k+1)} = b - Ax^{(k+1)}$ при $k \geq 0$ вычисляются по следующим формулам:

$$x^{(k+1)} = x^{(k)} + \alpha_k p^{(k)}, \quad \alpha_k = \frac{\|r^{(k)}\|_2^2}{(Ap^{(k)}, p^{(k)})}, \quad (6.67)$$

$$r^{(k+1)} = r^{(k)} - \alpha_k Ap^{(k)}, \quad (6.68)$$

$$p^{(k+1)} = r^{(k+1)} + \beta_k p^{(k)}, \quad \beta_k = \frac{\|r^{(k+1)}\|_2^2}{\|r^{(k)}\|_2^2}. \quad (6.69)$$

¹ В § 10.3 метод сопряженных градиентов рассматривается как метод минимизации квадратичной функции (6.63). Там же можно найти объяснение названия метода.

² Hestenes M.R., Stiefel E. L. Methods of conjugate gradients for solving linear systems // J. Res. Nat. Bur. Stand., Section B. 1952. V. 49. P. 409—436.

³ Reid J. K. On the method of conjugate gradients for the solution of large sparse systems of linear equations // In: J. K. Reid, editor. Large Sparse Sets of Linear Equations, Academic Press. 1971. pp. 231 — 254.

Теорема 6.11. Справедлива следующая оценка погрешности метода сопряженных градиентов:

$$\|x^{(n)} - \bar{x}\|_A \leq \frac{2\rho^n}{1 + \rho^{2n}} \|x^{(0)} - \bar{x}\|_A, \quad (6.70)$$

где $\rho = \frac{\sqrt{\text{cond}_2(A)} - 1}{\sqrt{\text{cond}_2(A)} + 1}$. Как следствие, справедлива оценка

$$\|x^{(n)} - \bar{x}\|_2 \leq \sqrt{\text{cond}_2(A)} \frac{2\rho^n}{1 + \rho^{2n}} \|x^{(0)} - \bar{x}\|_2. \blacksquare \quad (6.71)$$

Скорость сходимости метода сопряженных градиентов та же, что и у явного итерационного метода с чебышевским набором параметров, но в отличие от последнего для его реализации не нужна информация о значениях λ_{\min} и λ_{\max} .

Интересно, что метод сопряженных градиентов имеет все черты быстро сходящегося итерационного метода, но в то же время он является и прямым (точным) методом.

Теорема 6.12. Метод сопряженных градиентов позволяет найти точное решение системы $Ax = b$ с симметричной положительно определенной матрицей порядка t не более чем за t итераций. ■

Хотя теоретически одно из приближений $x^{(n)}$ с номером $n \leq t$ должно совпасть с точным решением \bar{x} , на практике при решении систем высокого порядка t этим свойством не удается воспользоваться. Во первых, получение точного решения невозможно из-за вычислительной погрешности. Во-вторых, как правило, для систем большого порядка t применение итерационного метода целесообразно только, если требуемое число итераций n много меньше t .

2. Сверхлинейная скорость сходимости. Классическая априорная оценка (6.70), казалось бы, указывает на то, что метод сопряженных градиентов обладает линейной скоростью сходимости. Однако эта оценка часто бывает чрезмерно пессимистичной и не объясняет тонких свойств сходимости, присущих методу сопряженных градиентов.

Для метода с линейной скоростью сходимости отношение $q_k = \frac{\|x^{(k+1)} - \bar{x}\|_A}{\|x^{(k)} - \bar{x}\|_A}$ практически не меняется с ростом k . Тем не менее,

опыт применения метода сопряженных градиентов показал, что это отношение обычно убывает с ростом k , и в ходе итерационного процесса наблюдается ускорение скорости сходимости метода.

Первое объяснение этого эффекта, получившего название *сверхлинейной сходимости*, на качественном уровне было дано в 1976 г.¹ В 1986 г. Вандерслюис и Вандерворт² дали его теоретическое обоснование. Оказалось, что метод сопряженных градиентов (в точной арифметике) через некоторое время начинает вести себя так, будто в матрице A исчезли крайние собственные значения $\lambda_1 = \lambda_{\min}(A)$ и $\lambda_n = \lambda_{\max}(A)$; затем — так, будто исчезли собственные значения λ_2 и λ_{n-1} и так далее. Таким образом, влияние крайних собственных значений оказывается только на первых итерациях и существенным оказывается распределение основной части собственных значений.

Исследование феномена сверхлинейной сходимости продолжается.

3. Метод сопряженных градиентов с предобусловливанием. Метод сопряженных градиентов сходится очень быстро, если $\text{cond}_2(A) \approx 1$. Напротив, метод сходится довольно медленно, если $\text{cond}_2(A) \gg 1$. Поэтому возникает естественное желание выполнить предобусловливание — заменить исходную систему уравнений $Ax = b$ эквивалентной системой

$$\tilde{A}y = \tilde{b} \quad (6.72)$$

с матрицей \tilde{A} , число обусловленности которой много меньше числа обусловленности матрицы A .

Пусть B — некоторая невырожденная матрица. Положим $y = Bx$ и умножим слева левую и правую части системы $Ax = b$ на $(B^T)^{-1}$. В результате получится система (6.72), в которой $\tilde{A} = (B^T)^{-1}AB^{-1}$, $\tilde{b} = (B^T)^{-1}b$.

Заметим, что построенная таким образом матрица \tilde{A} — симметричная и положительно определенная. Если $\text{cond}_2(\tilde{A}) \ll \text{cond}_2(A)$, то имеет смысл применить метод сопряженных градиентов (6.67)–(6.69) к преобразованной системе. Пусть $y^{(0)}$ — начальное приближение, $q^{(0)} = s^{(0)} = \tilde{b} - \tilde{A}y^{(0)}$, где $s^{(0)} = \tilde{b} - \tilde{A}y^{(0)}$ — начальное значение невязки. Прибли-

¹ Concus P., Golub G.H., O'Leary D.P. A generalized conjugate gradient method for the numerical solution of elliptic partial differential equations // In: J.R. Bunch and D.J. Rose, editors, Sparse Matrix Computations Academic Press, New York. 1976.

² Van der Sluis A., van der Vorst H.A. The rate of convergence of conjugate gradients // Numer. Math. 1986. V. 48. pp. 543–560.

жения $y^{(k)}$ к решению и невязки $s^{(k)} = \tilde{\mathbf{b}} - \tilde{\mathbf{A}}y^{(k)}$ при $k \geq 1$ вычисляются по формулам:

$$\mathbf{y}^{(k+1)} = \mathbf{y}^{(k)} + \alpha_k \mathbf{q}^{(k)}, \quad \alpha_k = \frac{\|s^{(k)}\|_2^2}{(\tilde{\mathbf{A}}\mathbf{q}^{(k)}, \mathbf{q}^{(k)})},$$

$$\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} - \alpha_k \tilde{\mathbf{A}}\mathbf{q}^{(k)},$$

$$\mathbf{q}^{(k+1)} = \mathbf{s}^{(k+1)} + \beta_k \mathbf{q}^{(k)}, \quad \beta_k = \frac{\|\mathbf{s}^{(k+1)}\|_2^2}{\|\mathbf{s}^{(k)}\|_2^2}.$$

Используя теперь формулы связи $\mathbf{y}^{(k)} = \mathbf{Bx}^{(k)}$, $\mathbf{s}^{(k)} = (\mathbf{B}^T)^{-1}\mathbf{r}^{(k)}$, $\mathbf{q}^{(k)} = \mathbf{Bp}^{(k)}$, приходим к методу сопряженных градиентов с предобусловливанием:

$$\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{Ax}^{(0)}, \quad \mathbf{B}^T \mathbf{s}^{(0)} = \mathbf{r}^{(0)}, \quad \mathbf{Bp}^{(0)} = \mathbf{s}^{(0)}, \quad (6.73)$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}, \quad \alpha_k = \frac{\|\mathbf{s}^{(k)}\|_2^2}{(\mathbf{Ap}^{(k)}, \mathbf{p}^{(k)})}, \quad (6.74)$$

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k \mathbf{Ap}^{(k)}, \quad (6.75)$$

$$\mathbf{B}^T \mathbf{s}^{(k+1)} = \mathbf{r}^{(k+1)}, \quad \mathbf{Bz}^{(k+1)} = \mathbf{s}^{(k+1)}, \quad (6.76)$$

$$\mathbf{p}^{(k+1)} = \mathbf{z}^{(k+1)} + \beta_k \mathbf{p}^{(k)}, \quad \beta_k = \frac{\|\mathbf{s}^{(k+1)}\|_2^2}{\|\mathbf{s}^{(k)}\|_2^2}. \quad (6.77)$$

Теорема 6.13. Для метода сопряженных градиентов с предобусловливанием (6.73)–(6.77) справедливы оценки погрешности (6.70),

$$(6.71) \quad c \rho = \frac{\sqrt{\operatorname{cond}_2(\tilde{\mathbf{A}})} - 1}{\sqrt{\operatorname{cond}_2(\tilde{\mathbf{A}})} + 1}. \blacksquare$$

Обратим внимание на то, что платой за ускорение сходимости метода стала необходимость решать на каждой итерации системы (6.76). Поэтому эффективность работы метода самым существенным образом зависит от того, как быстро могут решаться эти системы.

Нет общего рецепта, как построить предобусловливатель \mathbf{B} . Более того, его построение предполагает выполнение противоречивых свойств. С одной стороны, матрица \mathbf{B} должна быть такой, чтобы величина $\operatorname{cond}_2(\tilde{\mathbf{A}})$ была, по возможности, минимальной. В этом смысле идеальным был бы выбор $\mathbf{B} = \mathbf{L}^T$, где \mathbf{L} — невырожденная нижняя треугольная матрица, входящая в разложение Холецкого $\mathbf{A} = \mathbf{LL}^T$ (см. § 5.8).

При таком выборе $\tilde{A} = E$, $\text{cond}_2(\tilde{A}) = 1$ и метод сопряженных градиентов даст точный результат за один шаг. Однако, если мы уже имеем матрицу L и можем быстро решать системы вида $L^T y = b$, $Lx = y$, то в применении итерационного метода нет необходимости.

С другой стороны, необходимо, чтобы структура матрицы B была достаточно простой, чтобы быстро решать системы (6.76).

Ясно также, что вычисление предобусловливателя B не должно требовать большой вычислительной работы.

Одна из наиболее важных стратегий предобусловливания состоит в использовании неполного разложения Холецкого, в котором $B^T = \tilde{L}$ — аппроксимация матрицы L , входящей в разложение Холецкого. Простой метод вычисления \tilde{L} указан в [31].

Замечание. Свойство сверхлинейной сходимости (см. п. 2) подсказывает и другую стратегию предобусловливания. Предобусловливатель нужно выбирать не так, чтобы число обусловленности матрицы \tilde{A} было много меньше числа обусловленности матрицы A , а так, чтобы большая часть собственных значений матрицы \tilde{A} сконцентрировалась (кластеризовалась) вблизи единицы, то есть в интервале $(1 - \varepsilon_1, 1 + \varepsilon_2)$. Если вне этого интервала окажется лишь небольшой набор собственных значений, то после нескольких итераций метод сопряженных градиентов начнет вести себя так, как будто число обусловленности матрицы \tilde{A} не превышает $\frac{1 + \varepsilon_1}{1 - \varepsilon_2}$.

При этом возможно, что даже $\text{cond}_2(\tilde{A}) > \text{cond}_2(A)$.

§ 6.6. Дополнительные замечания

1. В этой главе были рассмотрены лишь основные итерационные методы решения систем линейных алгебраических уравнений. Вне нашего поля зрения остались многие известные и популярные методы. Среди них методы расщепления (неявные методы переменных направлений), блочные итерационные методы, методы подпространств Крылова, метод сопряженных невязок, обобщенный метод минимальных невязок (GMRES), методы MINRES, QMR, Bi-CGSTAB.

Эти методы изложены, например, в учебниках [9, 88, 100, 122, 2*, 5*, 16*] и в специальной литературе [24, 89, 110, 121, 123, 24*, 26*].

2. В гл. 10 мы еще раз вернемся к методам наискорейшего спуска и сопряженных градиентов в связи с задачей минимизации квадратичной функции (6.63). Обсуждение некоторых интересных особенностей метода сопряженных градиентов можно найти в [100, 5*], см. также [31].

3. В настоящее время наиболее глубоко развиты методы решения систем уравнений с симметричными положительно определенными матрицами. Однако часто встречаются задачи, в которых матрица A — несимметричная или симметричная, но не является положительно определенной. Ряд эффективных методов решения систем уравнений с несимметричными матрицами рассмотрен в [89, 110, 122, 123, 2*, 26*].

В принципе, всегда существует возможность симметризовать любую систему $Ax = b$ с невырожденной матрицей, т.е. свести ее к эквивалентной системе с симметричной положительно определенной матрицей. Для этого достаточно умножить обе части системы на матрицу A^T . В полученной таким образом системе

$$A^T A x = A^T b \quad (6.78)$$

матрица $\tilde{A} = A^T A$ обладает всеми нужными свойствами. Однако, как правило, это делать не рекомендуется. Дело в том, что при переходе от A к \tilde{A} может быть потеряно свойство разреженности. Как следствие, возрастают вычислительные затраты на одной итерации. Кроме того, как правило, существенно ухудшается обусловленность системы. Например, для матриц, близких к симметричным, $\text{cond}(\tilde{A}) \approx (\text{cond}(A))^2$. В силу этого даже самые быстрые методы решения систем с симметричными матрицами в применении к системе (6.78) могут стать очень незэффективными.

4. Как уже было отмечено, одно из важнейших достоинств итерационных методов заключается в возможности эффективного использования разреженности матрицы системы. Поясним сказанное на примере метода простой итерации (6.6).

В случае, когда матрица B — заполненная, для вычисления по формуле (6.6) требуется выполнить примерно $2m^2$ арифметических операций. Однако для разреженной матрицы с M ($M \ll m^2$) ненулевыми элементами требуется лишь примерно $2M$ арифметических операций (одно умножение и одно сложение на каждый ненулевой элемент). Таким образом, общее число операций составляет примерно $2Mn(\varepsilon)$, где $n(\varepsilon)$ — число итераций, необходимое для достижения заданной точности ε .

5. В случае, когда препятствием для быстрой сходимости итерационного метода является плохая обусловленность матрицы A , для ускорения сходимости систему $Ax = b$ заменяют эквивалентной системой $\tilde{A}y = \tilde{b}$ с матрицей вида $\tilde{A} = B_l^{-1} A B_r^{-1}$, которая имеет существенно меньшее число обусловленности. Этот подход, называемый *предобусловливанием*, уже обсуждался в § 6.4.

Как правило, матрицу \tilde{A} явным образом не вычисляют, используя в итерационном процессе наряду с исходной матрицей A невырожденные матрицы B_l и B_r , которые называют *левым* и *правым неявным предобусловливателями*. Неявными предобусловливатели называют потому, что их использование сводится к вычислению векторов вида $w = B_l^{-1}z$, $w = B_r^{-1}z$, т.е. к решению систем $B_l w = z$, $B_r w = z$. Если в преобразованной системе матрица \tilde{A} имеет вид $\tilde{A} = B_l A B_r$, то предобусловливатели называются *явными*, так как их применение сводится к явному умножению B_l и B_r , на соответствующие векторы. Возможны варианты $\tilde{A} = B_l^{-1} A B_r$, или $\tilde{A} = B_l A B_r^{-1}$, когда один из предобусловливателей является явным, а другой неявным. Часто один из предобусловливателей отсутствует; это не противоречит общей ситуации, а означает, что $B_l = E$ либо $B_r = E$.

При построении неявного метода с оптимальным параметром (6.51) теоретически система $Ax = b$ преобразовывалась к эквивалентной системе с матрицей $\tilde{A} = C^{-1}AC^{-1}$, где $C = \sqrt{B}$ играет роль левого и правого неявного предобусловливателей. Однако при численной реализации метода фактически система преобразуется к виду $B^{-1}Ax = B^{-1}b$, т.е. роль неявного левого предобусловливателя играет матрица B , а правый предобусловливателей отсутствует. Так же используются предобусловливатели в неявном методе с чебышевским набором параметров.

В методе сопряженных градиентов с предобусловливанием используется преобразование к виду (6.72), где $\tilde{A} = (B^T)^{-1}AB^{-1}$, т.е. оба предобусловливателя $B_l = B^T$ и $B_r = B$ — неявные. Важно то, что матрица $\tilde{A} = (B^T)^{-1}AB^{-1}$ — симметричная и положительно определенная, если исходная матрица A — симметричная и положительно определенная.

Обычно неявные предобусловливатели B_l и B_r выбираются так, чтобы системы уравнений с ними легко решались. В то же время желательно, чтобы в определенном смысле матрица \tilde{A} была как можно ближе к единичной матрице E , т.е. матрица $B_l B_r$ в некотором смысле должна приближать матрицу A .

6. Иногда, используя информацию о уже полученных приближениях к решению, удается построить уточненное приближение без существенных дополнительных вычислительных затрат.

Пусть, например, $x^{(n-2)}$, $x^{(n-1)}$ и $x^{(n)}$ — три последовательных приближения, полученных с помощью стационарного двухслойного итера-

ционного метода. Согласно δ^2 -процессу ускорения сходимости уточненное приближение $\tilde{x}^{(n)}$ имеет вид

$$\tilde{x}^{(n)} = x^{(n)} + \frac{x^{(n)} - x^{(n-1)}}{1 - \gamma_n}, \quad \gamma_n = \frac{(x^{(n-1)} - x^{(n-2)}, x^{(n)} - x^{(n-1)})}{\|x^{(n)} - x^{(n-1)}\|_2^2}.$$

Приближение $\tilde{x}^{(n)}$ может быть использовано далее как начальное для следующих итераций. Как критерий применимости δ^2 -процесса может быть использовано приближенное равенство

$$\frac{|(x^{(n-1)} - x^{(n-2)}, x^{(n)} - x^{(n-1)})|}{\|x^{(n-1)} - x^{(n-2)}\|_2 \|x^{(n)} - x^{(n-1)}\|_2} \approx 1.$$

Подробнее о δ^2 -процессе см. в [9].

Приближения $x^{(n-2)}$, $x^{(n-1)}$ и $x^{(n)}$ могут быть использованы и иначе. Будем искать уточненное приближение к решению в виде

$$\tilde{x}^{(n)} = x^{(n-2)} + \alpha(x^{(n-1)} - x^{(n-2)}) + \beta(x^{(n)} - x^{(n-2)}), \quad (6.75)$$

где параметры α и β определяются из условия минимума евклидовой нормы невязки $\tilde{r}^{(n)} = b - Ax^{(n)}$. Заметим, что

$$\tilde{r}^{(n)} = r^{(n-2)} + \alpha \Delta r^{(n-1)} + \beta \Delta r^{(n)},$$

где $r^{(n)} = b - Ax^{(n)}$, $\Delta r^{(n-1)} = r^{(n-1)} - r^{(n-2)}$, $\Delta r^{(n)} = r^{(n)} - r^{(n-2)}$.

Необходимое условие экстремума для функции

$$\Phi(\alpha, \beta) = \|r^{(n-2)} + \alpha \Delta r^{(n-1)} + \beta \Delta r^{(n)}\|_2^2.$$

приводит к системе уравнений

$$(\Delta r^{(n-1)}, \Delta r^{(n-1)})\alpha + (\Delta r^{(n-1)}, \Delta r^{(n)})\beta = -(r^{(n-2)}, \Delta r^{(n-1)}),$$

$$(\Delta r^{(n-1)}, \Delta r^{(n)})\alpha + (\Delta r^{(n)}, \Delta r^{(n)})\beta = -(r^{(n-2)}, \Delta r^{(n)})$$

относительно параметров α и β . Найденное приближение (6.75) таково, что

$$\|\tilde{r}^{(n)}\|_2 \leq \min\{\|r^{(n-2)}\|_2, \|r^{(n-1)}\|_2, \|r^{(n)}\|_2\}.$$

Дальнейшее развитие этой идеи приводит к обобщенному методу минимальных невязок (MINRES).

МЕТОДЫ ОТЫСКАНИЯ РЕШЕНИЙ СИСТЕМ НЕЛИНЕЙНЫХ УРАВНЕНИЙ

В этой главе рассматривается задача отыскания решений систем нелинейных уравнений, существенно более сложная, нежели задачи отыскания решения одного нелинейного уравнения или системы линейных алгебраических уравнений. Тем не менее достаточно подробное знакомство с содержанием гл. 4 и 6, а также § 5.1—5.3 позволяет увидеть соответствующие аналогии в постановках проблем и методах их решения для нелинейных систем.

Будем считать, что в множестве m -мерных векторов введена некоторая норма, порождающая соответствующую норму для квадратных матриц порядка m (см. § 5.2).

§ 7.1. Постановка задачи. Основные этапы решения

1. Постановка задачи. Задача отыскания решения системы нелинейных уравнений с m неизвестными вида

$$\begin{aligned} f_1(x_1, x_2, \dots, x_m) &= 0, \\ f_2(x_1, x_2, \dots, x_m) &= 0, \\ \dots & \\ f_m(x_1, x_2, \dots, x_m) &= 0 \end{aligned} \tag{7.1}$$

является существенно более сложной, чем рассмотренная в гл. 4 задача отыскания решения уравнения с одним неизвестным. Однако на практике она встречается значительно чаще, так как в реальных исследованиях интерес представляет, как правило, определение не одного, а нескольких параметров (нередко их число доходит до сотен и тысяч).

Найти точное решение системы, т.е. вектор $\bar{x} = (x_1, x_2, \dots, x_m)^T$, удовлетворяющий уравнениям (7.1), практически невозможно. В отличие от случая решения систем линейных алгебраических уравнений использование прямых методов здесь исключается. Единственно реальный путь решения системы (7.1) состоит в использовании итерационных методов для получения приближенного решения $x^* = (x_1^*, x_2^*, \dots, x_m^*)^T$, удовлетворяющего при заданном $\varepsilon > 0$ неравенству $\|x^* - \bar{x}\| < \varepsilon$.

Прежде чем перейти к изучению методов решения системы (7.1), подчеркнем важность понимания того факта, что эта задача может вообще не иметь решения, а если решения существуют, их число может быть произвольным. В общем случае весьма сложно выяснить, имеет ли система решения и сколько их.

Пример 7.1. Рассмотрим систему уравнений

$$4x_1^2 + x_2^2 = 4,$$

$$x_1 - x_2^2 + t = 0.$$

Здесь x_1, x_2 — неизвестные; t — параметр.

Первое уравнение задает на плоскости x_1Ox_2 эллипс, второе уравнение — параболу. Координаты точек пересечения этих кривых дают решения системы. Если значения параметра t изменяются от -2 до 2 , то возможны, например, следующие ситуации (рис. 7.1):

- а) $t = -2$ — решений нет;
- б) $t = -1$ — одно решение;
- в) $t = 0$ — два решения;
- г) $t = 1$ — три решения;
- д) $t = 2$ — четыре решения. ▲

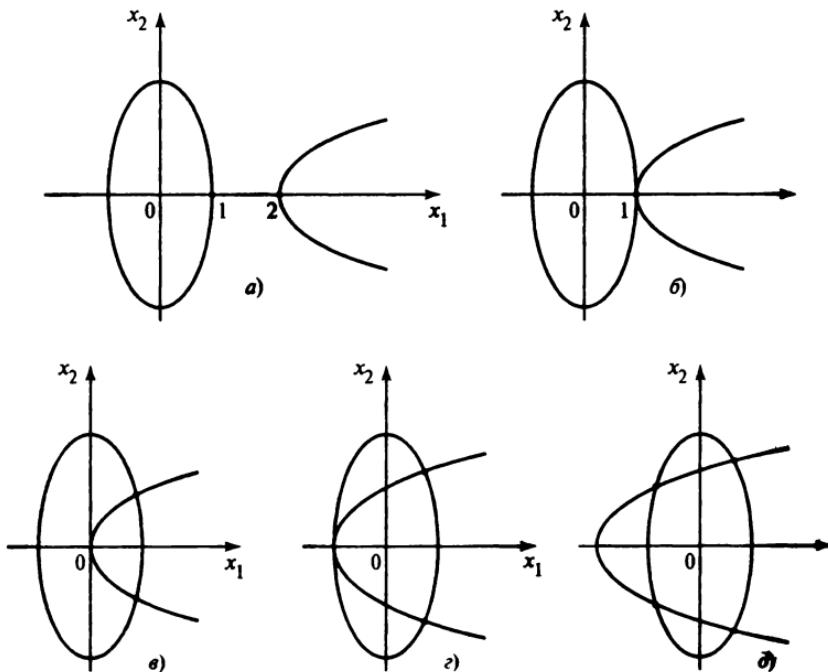


Рис. 7.1

Для дальнейшего удобно использовать сокращенную векторную форму записи систем. Наряду с вектором неизвестных $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$ рассмотрим вектор-функцию $\mathbf{f} = (f_1, f_2, \dots, f_m)^T$. В этих обозначениях система (7.1) примет вид

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}. \quad (7.2)$$

Будем считать функции $f_i(x)$ непрерывно дифференцируемыми в некоторой окрестности решения $\bar{\mathbf{x}}$. Введем для системы функций f_1, f_2, \dots, f_m матрицу Якоби

$$\mathbf{f}'(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_m} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_2(\mathbf{x})}{\partial x_m} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \frac{\partial f_m(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_m} \end{bmatrix}, \quad (7.3)$$

которая будет использована в дальнейшем.

2. Основные этапы решения. Как и для уравнения с одним неизвестным (см. гл. 4), отыскание решений начинают с этапа локализации. Для каждого из искомых решений $\bar{\mathbf{x}}$ определяют множество, содержащее только одно это решение и расположенное в достаточно малой его окрестности. Часто в качестве такого множества выступает параллелепипед или шар в m -мерном пространстве.

Иногда этап локализации не вызывает затруднений; соответствующие множества могут быть заданными, определяться из физических соображений, из смысла параметров x_i ; либо быть известными из опыта решений подобных задач. Однако чаще всего задача локализации (в особенности при больших m) представляет собой сложную проблему, от успешного решения которой в основном и зависит возможность вычисления решений системы (7.1). На этапе локализации особое значение приобретают квалификация исследователя, понимание им существа решаемой научной или инженерной проблемы, опыт решения этой или близких задач на компьютере. Во многих случаях полное решение задачи локализации невозможно и ее можно считать решенной удовлетворительно, если для $\bar{\mathbf{x}}$ удается найти хорошее начальное при-

ближение $x^{(0)}$. В простейших ситуациях (например, для системы двух уравнений с двумя неизвестными) могут быть использованы графические методы (см. пример 7.1).

На втором этапе для вычисления решения с заданной точностью ε используют один из итерационных методов решения нелинейных систем.

Будем считать, что определения § 4.1, связанные с характеризацией сходимости итерационных методов, остаются в силе, причем в неравенствах (4.5) и (4.6) знак модуля заменен на знак нормы, а δ -окрестность решения \bar{x} понимается как множество точек x , удовлетворяющих условию $\|x - \bar{x}\| < \delta$.

Пример 7.2. Произведем локализацию решений системы

$$\begin{aligned} x_1^3 + x_2^3 &= 8x_1x_2, \\ x_1 \ln x_2 &= x_2 \ln x_1. \end{aligned} \quad (7.4)$$

На плоскости x_10x_2 построим графики уравнений системы. График первого уравнения — это лист Декарта¹ (рис. 7.2, а). График второго уравнения состоит из луча — биссектрисы первого координатного угла и кривой, пересекающей эту биссектрису в точке $(e, e)^T$ (рис. 7.2, б).

Из рис. 7.3 видно, что эти кривые пересекаются в трех точках A, B, C , т.е. система имеет три решения. Так как оба графика симметричны относительно прямой $x_1 = x_2$, то координаты точки B равны и их легко вычислить:

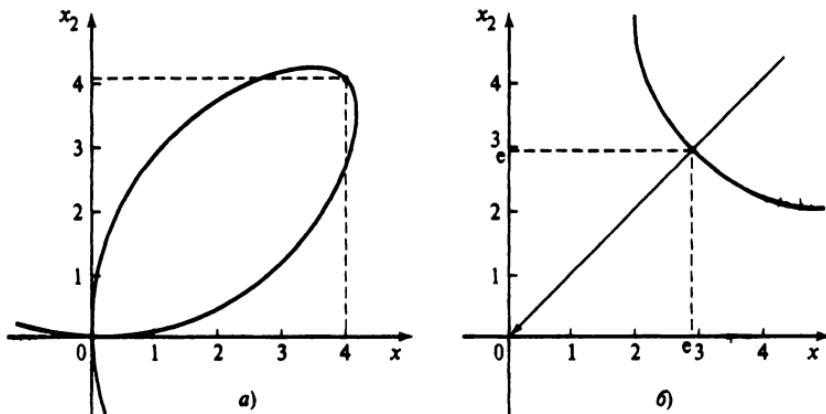


Рис. 7.2

¹ Рене Декарт (1596—1650) — французский философ, математик, физик, физиолог. Впервые ввел понятие переменной величины и функции. В аналитической геометрии создал метод прямоугольных координат.

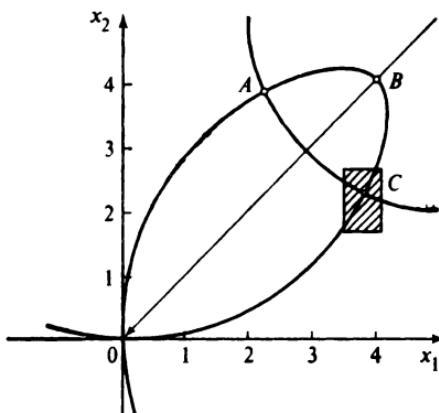


Рис. 7.3

$x_1 = 4, x_2 = 4$. В силу этой же симметрии достаточно определить только координаты \bar{x}_1, \bar{x}_2 точки C , так как точка A имеет координаты $x_1 = \bar{x}_2$ и $x_2 = \bar{x}_1$. Из рис. 7.3 замечаем, что точка C содержится в прямоугольнике $\Pi = \{(x, y) : 3.5 \leq x_1 \leq 4, 1.5 \leq x_2 \leq 2.5\}$ и $\bar{x}_1 \approx 3.8, \bar{x}_2 \approx 2.5$. \blacktriangle

Подчеркнем, что только по виду уравнений системы (7.4) без использования графического метода установить число решений и найти приближения к ним было бы весьма трудно. К сожалению, при числе уравнений $m > 3$ геометрические иллюстрации теряют свою эффективность.

Замечание. Иногда удается понизить порядок m системы, выразив одно или несколько неизвестных из одних уравнений системы и подставив соответствующие выражения в другие уравнения.

Пример 7.3. Система уравнений

$$x^3 + y^3 = 8xy,$$

$$x \ln y = \ln x$$

сводится к одному нелинейному уравнению $x^3 + x^{3/x} = 8x^1 + 1/x$ после того, как из второго уравнения выражается $y = x^{1/x}$. \blacktriangle

3. Корректность и обусловленность задачи. Будем считать, что система (7.1) имеет решение \bar{x} , причем в некоторой окрестности этого решения матрица Якоби $f'(x)$ невырождена. Выполнение последнего условия гарантирует, что в указанной окрестности нет других решений системы (7.1). Случай, когда в точке \bar{x} матрица $f'(x)$ вырождена,

является существенно более трудным и нами рассматриваться не будет. В одномерном случае первая ситуация отвечает наличию простого корня уравнения $f(x) = 0$, а вторая — кратного корня.

В § 4.2 было показано, что погрешность вычисления функции f приводит к образованию вокруг корня уравнения $f(x) = 0$ интервала неопределенности, внутри которого невозможно определить, какая из точек является решением уравнения.

Аналогично, погрешности в вычислении вектор-функции f приводят к появлению *области неопределенности* D , содержащей решение \bar{x} системы (7.1) такой, что для всех $x \in D$ векторное уравнение $f(x) = 0$ удовлетворяется с точностью до погрешности. Область D может иметь довольно сложную геометрическую структуру (рис. 7.4). Мы удовлетворимся только лишь оценкой радиуса $\bar{\varepsilon}$ этой области.

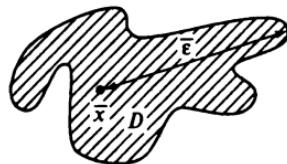


Рис. 7.4

Предположим, что для близких к \bar{x} значений x вычисляемые значения $f^*(x)$ удовлетворяют неравенству $\|f(x) - f^*(x)\| \leq \bar{\Delta}(f^*)$. Тогда $\bar{\varepsilon}$ можно приближенно оценить с помощью неравенства $\bar{\varepsilon} \leq \|(f'(\bar{x}))^{-1}\| \bar{\Delta}(f^*)$. Таким образом, в рассматриваемой задаче роль абсолютного числа обусловленности играет норма матрицы, обратной матрице Якоби $f'(\bar{x})$.

§ 7.2. Метод простой итерации

1. Описание метода. Предположим, что требуется найти решение $\bar{x} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m)^T$ системы (7.1) с заданной точностью $\varepsilon > 0$. Преобразуем систему (7.1) к следующему эквивалентному виду (к виду, удобному для итераций):

$$\begin{aligned} x_1 &= \varphi_1(x_1, x_2, \dots, x_m), \\ x_2 &= \varphi_2(x_1, x_2, \dots, x_m), \\ &\dots \\ x_m &= \varphi_m(x_1, x_2, \dots, x_m). \end{aligned} \tag{7.5}$$

Если ввести вектор-функцию $\varphi = (\varphi_1, \varphi_2, \dots, \varphi_m)^T$, то система (7.5) запишется так:

$$\mathbf{x} = \varphi(\mathbf{x}). \quad (7.6)$$

Пусть начальное приближение $\mathbf{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_m^{(0)})^T$ задано.

Подставляя его в правую часть системы (7.6), получаем $\mathbf{x}^{(1)} = \varphi(\mathbf{x}^{(0)})$. Подставляя $\mathbf{x}^{(1)}$ в правую часть (7.6), находим $\mathbf{x}^{(2)} = \varphi(\mathbf{x}^{(1)})$ и т.д. Продолжая вычисления по формулам

$$\mathbf{x}^{(k+1)} = \varphi(\mathbf{x}^{(k)}), \quad k \geq 0, \quad (7.7)$$

получаем последовательность $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}, \dots$ приближений к решению $\bar{\mathbf{x}}$.

Запись (7.7) означает, что очередное приближение $\mathbf{x}^{(k+1)}$ вычисляется через предыдущее приближение $\mathbf{x}^{(k)}$ следующим образом:

$$\begin{aligned} x_1^{(k+1)} &= \varphi_1(x_1^{(k)}, x_2^{(k)}, \dots, x_m^{(k)}), \\ x_2^{(k+1)} &= \varphi_2(x_1^{(k)}, x_2^{(k)}, \dots, x_m^{(k)}), \\ &\dots \\ x_m^{(k+1)} &= \varphi_m(x_1^{(k)}, x_2^{(k)}, \dots, x_m^{(k)}). \end{aligned}$$

Отметим существенную аналогию с методами простой итерации для решения одного нелинейного уравнения (см. гл. 4) и системы линейных алгебраических уравнений (см. гл. 6).

2. Сходимость метода. Пусть $\varphi'(\mathbf{x})$ — матрица Якоби, отвечающая вектор-функции $\varphi(\mathbf{x})$ (см. § 7.1). Сформулируем теорему о сходимости метода простых итераций, являющуюся аналогом теорем 4.2 и 6.1.

Теорема 7.1. Пусть в некоторой σ -окрестности решения $\bar{\mathbf{x}}$ функции $\varphi_i(\mathbf{x})$ ($i = 1, 2, \dots, m$) дифференцируемы и выполнено неравенство

$$\|\varphi'(\mathbf{x})\| \leq q, \quad (7.8)$$

где $0 \leq q < 1$, q — постоянная.

Тогда независимо от выбора начального приближения $\mathbf{x}^{(0)}$ из указанной σ -окрестности корня итерационная последовательность не выходит из этой окрестности, метод сходится со скоростью геометрической прогрессии и справедлива следующая оценка погрешности:

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| \leq q^n \|\mathbf{x}^{(0)} - \bar{\mathbf{x}}\|. \blacksquare \quad (7.9)$$

Замечание 1. Условие (7.8), грубо говоря, означает, что в окрестности решения производные $\partial\phi_i/\partial x_j$ для всех i и j должны быть «достаточно малы по модулю». Иными словами, систему (7.1) следует преобразовать к такому виду (7.5), чтобы функции ϕ , слабо менялись при изменении аргументов, т.е. были «почти постоянными». Каких-либо общих рецептов, как это следует делать, в общем случае нет.

Замечание 2. В условиях теоремы 7.1 верна апостериорная оценка погрешности

$$\|\mathbf{x}^{(n)} - \bar{\mathbf{x}}\| \leq \frac{q}{1-q} \|\mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}\|,$$

которая удобна для формулирования критерия окончания итераций, если известна величина q . Однако найти значение q , удовлетворяющее неравенству (7.8) для всех x из некоторой σ -окрестности корня, далеко не просто. В ряде случаев при наличии достаточно хорошего начального приближения $\mathbf{x}^{(0)}$ можно, считая, что $q \approx q_0 = \|\phi'(\mathbf{x}^{(0)})\|$, использовать следующий практический критерий окончания итерационного процесса:

$$\|\mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}\| \lesssim \varepsilon_1 = \frac{1-q_0}{q_0} \varepsilon. \quad (7.10)$$

Пример 7.4. Используя метод простой итерации, найдем с точностью $\varepsilon = 10^{-3}$ решение \bar{x}_1, \bar{x}_2 системы (7.4).

Приведем систему к виду, удобному для итераций:

$$x_1 = \sqrt[3]{8x_1x_2 - x_2^3},$$

$$x_2 = x_2 + \frac{x_2}{\ln x_2} - \frac{x_1}{\ln x_1};$$

здесь $\phi_1(x_1, x_2) = \sqrt[3]{8x_1x_2 - x_2^3}$, $\phi_2(x_1, x_2) = x_2 + \frac{x_2}{\ln x_2} - \frac{x_1}{\ln x_1}$.

Проверим выполнение условия сходимости вблизи точки С (см. рис. 7.3). Вычислим матрицу Якоби

$$\phi'(\mathbf{x}_1, \mathbf{x}_2) = \begin{bmatrix} \frac{\partial \phi_1}{\partial x_1} & \frac{\partial \phi_1}{\partial x_2} \\ \frac{\partial \phi_2}{\partial x_1} & \frac{\partial \phi_2}{\partial x_2} \end{bmatrix} = \begin{bmatrix} \frac{8x_2}{3(8x_1x_2 - x_2^3)^{2/3}} & \frac{8x_1 - 3x_2^2}{3(8x_1x_2 - x_2^3)^{2/3}} \\ \frac{1}{\ln^2 x_1} - \frac{1}{\ln x_1} & 1 + \frac{1}{\ln x_2} - \frac{1}{\ln^2 x_2} \end{bmatrix}.$$

Так как $\bar{x}_1 \approx 3.8$ и $\bar{x}_2 \approx 2$, то для $x \approx \bar{x} \approx (3.8, 2)$ имеем

$$\varphi'(x_1, x_2) \approx \varphi'(3.8, 2) \approx \begin{bmatrix} 0.379 & 0.436 \\ -0.188 & 0.361 \end{bmatrix}$$

Тогда $\|\varphi'(x_1, x_2)\|_\infty \approx \|\varphi'(3.8, 2)\|_\infty \approx 0.379 + 0.436 = 0.815$. Следовательно, метод простой итерации

$$\begin{aligned} x_1^{(k+1)} &= \sqrt[3]{8x_1^{(k)}x_2^{(k)} - (x_2^{(k)})^3}, \\ x_2^{(k+1)} &= x_2^{(k)} + \frac{x_2^{(k)}}{\ln x_2^{(k)}} - \frac{x_1^{(k)}}{\ln x_1^{(k)}} \end{aligned} \quad (7.11)$$

будет сходиться со скоростью геометрической прогрессии, знаменатель которой примерно равен 0.815, если начальные приближения брать в достаточно малой окрестности решения

Возьмем $x_1^{(0)} = 3.8$, $x_2^{(0)} = 2$ и будем вести итерации по формулам (7.11), используя критерий окончания (7.10), в котором $\varepsilon = 10^{-3}$, $q_0 = 0.815$. Результаты вычислений с шестью знаками мантиссы приведены в табл. 7.1.

Таблица 7.1

k	0	1	2	3	4
$x_1^{(k)}$	3 80000	3 75155	3 74960	3 75892	3 76720
$x_2^{(k)}$	2 00000	2 03895	2 06347	2 07498	2 07883

Продолжение табл. 7.1

k	5	6	7	8	9
$x_1^{(k)}$	3 77198	3 77399	3 77450	3 77440	3 77418
$x_2^{(k)}$	2 07920	2 07850	2 07778	2 07732	2 07712

При $k = 9$ критерий окончания выполняется и можно положить $\bar{x}_1 = 3.774 \pm 0.001$, $\bar{x}_2 = 2.077 \pm 0.001$. ▲

3. Влияние погрешности вычислений. В силу наличия погрешностей при вычислении на компьютере получается не последовательность $x^{(k)}$, удовлетворяющая равенству (7.7), а последовательность $\tilde{x}^{(k)}$, удовлетворяющая равенству

$$\tilde{x}^{(k+1)} = \varphi^*(\tilde{x}^{(k)}). \quad (7.12)$$

Будем считать, что абсолютная погрешность вычисляемых значений $\varphi^*(x)$ вектор-функции φ мала и что $\|\varphi(x) - \varphi^*(x)\| \leq \bar{\Delta}(\varphi^*)$. Наличие погрешности вычисления φ приводит к появлению области неопределенности решения \bar{x} , радиус $\bar{\varepsilon}$ которой можно приближенно оценить, воспользовавшись неравенством $\bar{\varepsilon} \lesssim \bar{\varepsilon}^* = \bar{\Delta}(\varphi^*)/(1-q)$ в том случае, если $\|\varphi'(x)\| \leq q$.

Сформулируем следующий результат, являющийся аналогом теорем 4.5 и 6.3.

Теорема 7.2. Пусть выполнены условия теоремы 7.1 и для всех x из σ -окрестности решения \bar{x} выполнено неравенство $\|\varphi(x) - \varphi^*(x)\| \leq \bar{\Delta}(\varphi^*)$. Предположим также, что $\bar{\varepsilon}^* = \frac{\bar{\Delta}(\varphi^*)}{1-q} < \sigma$

(т.е. величина $\bar{\Delta}(\varphi^*)$ достаточно мала). Тогда если вычисления по формулам (7.7) и (7.12) начинаются с одного начального приближения $x^{(0)} = \tilde{x}^{(0)}$, принадлежащего σ_1 -окрестности решения \bar{x} (где $\sigma_1 = \min\{\sigma, (\sigma - \bar{\varepsilon}^*)/q\}$), то последовательность $\tilde{x}^{(n)}$ не выходит за пределы σ -окрестности решения \bar{x} и для всех $n \geq 1$ справедливы оценки

$$\begin{aligned}\|x^{(n)} - \tilde{x}^{(n)}\| &\leq \bar{\varepsilon}^*, \\ \|\tilde{x}^{(n)} - \bar{x}\| &\leq q^n \|\tilde{x}^{(0)} - \bar{x}\| + \bar{\varepsilon}^*, \\ \|\tilde{x}^{(n)} - \bar{x}\| &\leq \frac{q}{1-q} \|\tilde{x}^{(n)} - \tilde{x}^{(n-1)}\| + \bar{\varepsilon}^*. \blacksquare\end{aligned}$$

4. Модификации метода простой итерации. В некоторых случаях для ускорения сходимости полезно использовать следующий аналог метода Зейделя (см. § 6.2):

$$\begin{aligned}x_1^{(k+1)} &= \varphi_1(x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, \dots, x_m^{(k)}), \\ x_2^{(k+1)} &= \varphi_2(x_1^{(k+1)}, x_2^{(k)}, x_3^{(k)}, \dots, x_m^{(k)}), \\ x_3^{(k+1)} &= \varphi_3(x_1^{(k+1)}, x_2^{(k+1)}, x_3^{(k)}, \dots, x_m^{(k)}), \\ &\dots \\ x_m^{(k+1)} &= \varphi_m(x_1^{(k+1)}, x_2^{(k+1)}, x_3^{(k+1)}, \dots, x_m^{(k)}).\end{aligned}$$

Другой вариант метода Зейделя состоит в следующем: i -я компонента решения ($i = 1, 2, \dots, m$) на $(k + 1)$ -й итерации определяется как решение нелинейного уравнения

$$F_i(x_i^{(k+1)}) = 0,$$

где $F_i(x_i) = f_i(x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_i, x_{i+1}^{(k)}, \dots, x_m^{(k)})$.

Преимущества этого метода¹ состоят в возможности использования методов решения уравнений с одним неизвестным и в отсутствии необходимости приведения системы уравнений к виду, удобному для итераций. Указанный метод сходится, если для матрицы Якоби $f'(x)$ выполнены условия строчного диагонального преобладания.

Иногда существенный выигрыш дает использование *нелинейного SOR-метода*, являющегося аналогом метода релаксации (см. § 6.3). В нем после вычисления очередной i -й компоненты $(k + 1)$ -го приближения по формуле метода Зейделя

$$\hat{x}_i^{(k+1)} = \varphi_i(x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_i^{(k)}, \dots, x_m^{(k)})$$

приближение $x_i^{(k+1)}$ вычисляют по формуле

$$x_i^{(k+1)} = \hat{x}_i^{(k+1)} + (\omega - 1)(\hat{x}_i^{(k+1)} - x_i^{(k)}),$$

где ω — параметр релаксации.

§ 7.3. Метод Ньютона для решения систем нелинейных уравнений

1. Описание метода. Обобщим метод Ньютона, изложенный в § 4.6 для решения одного нелинейного уравнения, на решение систем нелинейных уравнений (7.1). При этом будем исходить из трактовки метода Ньютона как метода линеаризации.

Предположим, что исходя из начального приближения $x^{(0)}$ к решению \bar{x} построены приближения $x^{(1)}, x^{(2)}, \dots, x^{(n)}$. Заменим в системе (7.1) каждую из функций f_i ($i = 1, 2, \dots, m$) главной линейной частью ее разложения по формуле Тейлора в точке $x^{(n)}$:

$$f_i(\mathbf{x}) \approx f_i(\mathbf{x}^{(n)}) + \sum_{j=1}^m \frac{\partial f_i(\mathbf{x}^{(n)})}{\partial x_j} (x_j - x_j^{(n)}).$$

¹ Этот метод иногда называют *нелинейным методом Зейделя*.

В результате придем к системе линейных алгебраических уравнений

$$f_1(\mathbf{x}^{(n)}) + \sum_{j=1}^m \frac{\partial f_1(\mathbf{x}^{(n)})}{\partial x_j} (x_j - x_j^{(n)}) = 0,$$

$$f_2(\mathbf{x}^{(n)}) + \sum_{j=1}^m \frac{\partial f_2(\mathbf{x}^{(n)})}{\partial x_j} (x_j - x_j^{(n)}) = 0,$$

.....

$$f_m(\mathbf{x}^{(n)}) + \sum_{j=1}^m \frac{\partial f_m(\mathbf{x}^{(n)})}{\partial x_j} (x_j - x_j^{(n)}) = 0,$$

имеющей в матричной форме записи следующий вид:

$$\mathbf{f}(\mathbf{x}^{(n)}) + \mathbf{f}'(\mathbf{x}^{(n)})(\mathbf{x} - \mathbf{x}^{(n)}) = 0; \quad (7.13)$$

здесь \mathbf{f}' — матрица Якоби (7.3).

Предположим, что матрица $\mathbf{f}'(\mathbf{x}^{(n)})$ невырожденная, т.е. существует обратная матрица $(\mathbf{f}'(\mathbf{x}^{(n)}))^{-1}$. Тогда система (7.13) имеет единственное решение, которое и принимается за очередное приближение $\mathbf{x}^{(n+1)}$ к решению $\bar{\mathbf{x}}$. Таким образом, приближение $\mathbf{x}^{(n+1)}$ удовлетворяет равенству

$$\mathbf{f}(\mathbf{x}^{(n)}) + \mathbf{f}'(\mathbf{x}^{(n)})(\mathbf{x}^{(n+1)} - \mathbf{x}^{(n)}) = 0, \quad (7.14)$$

выражая из которого $\mathbf{x}^{(n+1)}$, выводим итерационную формулу метода Ньютона:

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - (\mathbf{f}'(\mathbf{x}^{(n)}))^{-1} \mathbf{f}(\mathbf{x}^{(n)}). \quad (7.15)$$

Замечание. Формула (7.15) предполагает использование трудоемкой операции обращения матрицы (см. гл. 5), поэтому непосредственное ее использование для вычисления $\mathbf{x}^{(n+1)}$ в большинстве случаев нецелесообразно. Обычно вместо этого решают эквивалентную системе (7.14) систему линейных алгебраических уравнений

$$\mathbf{f}'(\mathbf{x}^{(n)}) \Delta \mathbf{x}^{(n+1)} = -\mathbf{f}(\mathbf{x}^{(n)}) \quad (7.16)$$

относительно поправки $\Delta \mathbf{x}^{(n+1)} = \mathbf{x}^{(n+1)} - \mathbf{x}^{(n)}$. Затем полагают

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + \Delta \mathbf{x}^{(n+1)}. \quad (7.17)$$

2. Сходимость метода. Сформулируем основную теорему о сходимости метода Ньютона.

Теорема 7.3. Пусть в некоторой окрестности решения $\bar{\mathbf{x}}$ системы (7.1) функции f_i ($i = 1, 2, \dots, m$) дважды непрерывно дифференцируемы и матрица $\mathbf{f}'(\bar{\mathbf{x}})$ невырождена. Тогда найдется такая малая δ -окрест-

ность решения \bar{x} , что при произвольном выборе начального приближения $x^{(0)}$ из этой окрестности итерационная последовательность метода Ньютона не выходит за пределы окрестности и справедлива оценка:

$$\|x^{(n+1)} - \bar{x}\| \leq \delta^{-1} \|x^{(n)} - \bar{x}\|^2, \quad n \geq 0. \blacksquare$$

Эта оценка означает, что метод сходится с квадратичной скоростью.

Квадратичная скорость сходимости метода Ньютона позволяет использовать простой практический критерий окончания итераций

$$\|x^{(n)} - x^{(n-1)}\| < \varepsilon. \quad (7.18)$$

Пример 7.5. Используя метод Ньютона, найдем с точностью $\varepsilon = 10^{-4}$ решение \bar{x}_1, \bar{x}_2 системы (7.4).

Возьмем $x_1^{(0)} = 3.8, x_2^{(0)} = 2$ и будем вести вычисления по формулам (7.16), (7.17), в которых

$$f(x) = \begin{bmatrix} x_1^3 + x_2^3 - 8x_1x_2 \\ x_1 \ln x_2 - x_2 \ln x_1 \end{bmatrix}, \quad f'(x) = \begin{bmatrix} 3x_1^2 - 8x_2 & 3x_2^2 - 8x_1 \\ \ln x_2 - \frac{x_2}{x_1} & \frac{x_1}{x_2} - \ln x_1 \end{bmatrix}.$$

Результаты вычислений с шестью знаками мантиссы приведены в табл. 7.2.

Таблица 7.2

k	0	1	2	3
$x_1^{(k)}$	3.80000	3.77258	3.77388	3.77389
$x_2^{(k)}$	2.00000	2.07189	2.07708	2.07710

При $k = 3$ критерий окончания итераций $\|x^{(k)} - x^{(k-1)}\|_\infty < \varepsilon = 10^{-4}$ выполняется и можно положить $\bar{x}_1 = 3.7739 \pm 0.0001, \bar{x}_2 = 2.0771 \pm 0.0001$. \blacktriangle

3. Трудности использования метода Ньютона. Изложенные в § 4.6 трудности использования метода Ньютона не только сохраняются при применении его к решению систем нелинейных уравнений, но и усугубляются. Во-первых, возникает проблема вычисления на каждой итерации матрицы $f'(x)$ из m^2 частных производных, что само по себе может оказаться весьма сложным делом. Во-вторых, обостряется про-

блема нахождения хорошего начального приближения. Ее решить в многомерном случае гораздо труднее, чем в одномерном.

4. Модифицированный метод Ньютона. Модифицированный метод Ньютона призван существенно ослабить требования к начальному приближению. Опишем одну его итерацию.

Сначала, как и в классическом варианте метода, решают систему

$$f'(x^{(n)})\Delta x^{(n+1)} = -f(x^{(n)})$$

относительно вектора $\Delta x^{(n+1)}$.

Затем находят число α_n , решая задачу одномерной минимизации

$$F_n(\alpha_n) = \min_{\alpha} F_n(\alpha); \quad F_n(\alpha) = \|f(x^{(n)} + \alpha \Delta x^{(n+1)})\|.$$

Следующее приближение вычисляют по формуле

$$x^{(n+1)} = x^{(n)} + \alpha_n \Delta x^{(n+1)}.$$

Подробнее с модифицированным методом Ньютона можно ознакомиться в [104].

5. Влияние погрешности вычислений. Пусть f^* и $(f')^*$ — вычисляемые на компьютере приближенные значения вектор-функции f и матрицы Якоби f' . Пусть для решения системы линейных алгебраических уравнений (7.16) используется схема частичного выбора (см. § 5.5). Будем считать, что матрица f' достаточно хорошо обусловлена ($\text{cond}(f')\epsilon_m \ll 1$) и вычисляется не слишком грубо ($\|f' - (f')^*\| \ll \|f'\|$). Тогда при выборе начального приближения из малой окрестности решения метод Ньютона является устойчивым и дает возможность найти решение с гарантированной точностью $\epsilon \geq \bar{\epsilon}^* = \|(f'(\bar{x}))^{-1}\| \|\bar{\Delta}(f^*)\|$.

§ 7.4. Модификации метода Ньютона

Если оценивать качество метода Ньютона только по числу необходимых итераций, то следовало бы сделать вывод о том, что этот метод стоит применять всегда, когда он сходится. На практике для достижения разумной точности ϵ при выборе достаточно хорошего начального приближения $x^{(0)}$ требуется, как правило, три-пять итераций.

Однако при оценке общей трудоемкости метода следует учитывать, что на каждой итерации требуется выполнение следующей дополнительной работы:

- 1) вычисление m компонент вектора $f(x^{(n)})$;
- 2) вычисление m^2 компонент матрицы Якоби $f'(x^{(n)})$;
- 3) решение системы линейных алгебраических уравнений (7.16).

Существует большое число модификаций метода Ньютона, позволяющих в тех или иных ситуациях снизить его трудоемкость либо избежать необходимости вычисления производных. Рассмотрим кратко некоторые из таких модификаций, с учетом того, что их одномерные аналоги были изучены более подробно в § 4.7.

1. Упрощенный метод Ньютона. Заменим в расчетных формулах метода Ньютона (7.16), (7.17) матрицу $f'(x^{(n)})$, зависящую от n , постоянной матрицей $A = f'(x^{(0)})$. В результате получим расчетные формулы *упрощенного метода Ньютона*:

$$A\Delta x^{(n+1)} = -f(x^{(n)}), \quad (7.19)$$

$$x^{(n+1)} = x^{(n)} + \Delta x^{(n+1)}. \quad (7.20)$$

Можно показать, что этот метод сходится со скоростью геометрической прогрессии, если начальное приближение $x^{(0)}$ выбрано достаточно близким к решению \bar{x} , причем знаменатель прогрессии q тем меньше, чем ближе $x^{(0)}$ к \bar{x} .

По сравнению с методом Ньютона число итераций, необходимое для достижения заданной точности ϵ , существенно возрастает. Тем не менее общие вычислительные затраты могут оказаться меньше. Причины этого состоят в следующем. Во-первых, вычисление матрицы Якоби производится здесь только один раз; во-вторых, как нетрудно видеть, при использовании упрощенного метода Ньютона (7.19), (7.20) многократно решается система линейных уравнений с фиксированной матрицей A и различными правыми частями. Это означает, что при решении систем (7.19) методом Гаусса возможно применение *LU*-разложения матрицы A , которое резко уменьшает число операций, необходимых для вычисления $\Delta x^{(n+1)}$ (см. гл. 5).

2. Использование формул численного дифференцирования. Довольно часто вычисление производных $\partial f_i / \partial x_j$, являющихся элементами матрицы f' , затруднено или вообще невозможно. В такой ситуации для приближенного вычисления производных можно попытаться использовать формулы численного дифференцирования (см. гл. 12). Например, можно использовать следующую конечно-разностную аппроксимацию производной:

$$\frac{\partial f_i}{\partial x_j}(x^{(n)}) \approx J_{ij}^{(n)} =$$

$$= \frac{1}{h_j^{(n)}} (f_i(x_1^{(n)}, \dots, x_{j-1}^{(n)}, x_j^{(n)} + h_j^{(n)}, x_{j+1}^{(n)}, \dots, x_m^{(n)}) - f_i(x_1^{(n)}, \dots, x_m^{(n)})). \quad (7.21)$$

Параметры $h_j^{(n)} \neq 0$ ($j = 1, 2, \dots, m$) — это *конечно-разностные шаги*.

Если в расчетных формулах метода Ньютона (7.16), (7.17) заменить матрицу $f'(x^{(n)})$ аппроксимирующей ее матрицей $J^{(n)}$ с элементами $J_{ij}^{(n)}$, то получим следующий итерационный метод:

$$J^{(n)} \Delta x^{(n+1)} = -f(x^{(n)}), \quad (7.22)$$

$$x^{(n+1)} = x^{(n)} + \Delta x^{(n+1)}. \quad (7.23)$$

В простейшем варианте этого метода шаги h_j ($j = 1, 2, \dots, m$) не зависят от n . Отметим, что выбор величины шагов представляет собой не очень простую задачу. С одной стороны, они должны быть достаточно малыми, чтобы матрица $J^{(n)}$ хорошо приближала матрицу $f'(x^{(n)})$, с другой, они не могут быть очень малы, так как в этом случае влияние погрешностей вычисления функций f_j на погрешность формулы (7.21) численного дифференцирования становится катастрофическим (выполняется вычитание близких приближенных чисел).

Следующие три метода можно рассматривать как варианты метода (7.22), (7.23), в которых реализованы специальные подходы к вычислению вектора $h^{(n)} = (h_1^{(n)}, h_2^{(n)}, \dots, h_m^{(n)})^T$. Для того чтобы приведенные ниже рассуждения были формально корректными, в формуле (7.21) положим $J_{ij}^{(n)} = \frac{\partial f_i}{\partial x_j}(x^{(n)})$, если оказалось, что $h_j^{(n)} = 0$.

3. Метод ложного положения. Пусть c — фиксированный вектор. Положим $h^{(n)} = c - x^{(n)}$. Тогда формулы (7.21)–(7.23) определяют *метод ложного положения*, обладающий линейной скоростью сходимости, если вектор c и начальное приближение $x^{(0)}$ выбраны достаточно близко к решению.

4. Метод секущих. В одном из наиболее популярных своих вариантов *метод секущих* можно рассматривать как метод (7.21)–(7.23), где $h^{(n)} = x^{(n-1)} - x^{(n)}$. Метод секущих является двухшаговым: для вычисления очередного приближения $x^{(n+1)}$ используют два предыдущих приближения $x^{(n)}$ и $x^{(n-1)}$. Для того чтобы начать вычисления, необходимо задать два начальных приближения $x^{(0)}$ и $x^{(1)}$.

При удачном выборе начальных приближений $x^{(0)}$ и $x^{(1)}$ метод секущих сходится со сверхлинейной скоростью с порядком сходимости $p = (\sqrt{5} + 1)/2$.

5. Метод Стеффенсена. Вычисления по *методу Стеффенсена* производят по формулам (7.21)–(7.23), где $\mathbf{h}^{(n)} = \mathbf{f}(\mathbf{x}^{(n)})$. Замечательно то, что хотя этот метод не требует вычисления производных и в отличие от метода секущих является одношаговым, он, как и метод Ньютона, обладает свойством квадратичной сходимости. Правда, как и в методе Ньютона, его применение затруднено необходимостью выбора хорошего начального приближения. По-видимому, для решения нелинейных систем вида (7.1) метод Стеффенсена чаще окажется лучшим выбором, чем метод секущих или метод ложного положения.

Как и в одномерном случае (см. § 4.7), следует отметить, что методы секущих и Стеффенсена теряют устойчивость вблизи решения (фактически это происходит при попадании приближения $\mathbf{x}^{(n)}$ в область неопределенности решения $\bar{\mathbf{x}}$). Поэтому при использовании этих методов важно вовремя прекратить выполнение итераций.

6. Метод Бройдена. Существует класс *квазиньютоновских методов*, обобщающий на случай систем уравнений метод секущих. Среди этих методов заслуженной популярностью пользуется *метод Бройдена*.

В этом методе наряду с последовательностью приближений $\mathbf{x}^{(k)}$ к решению вычисляется последовательность приближений $\mathbf{B}^{(k)}$ к матрице Якоби $\mathbf{f}'(\bar{\mathbf{x}})$. Пусть заданы начальные приближения $\mathbf{x}^{(0)} \approx \bar{\mathbf{x}}$ и $\mathbf{B}^{(0)} \approx \mathbf{f}'(\bar{\mathbf{x}})$. Следующие приближения при $k = 0, 1, \dots$, вычисляются по формулам

$$\begin{aligned}\mathbf{B}^{(k)} \Delta \mathbf{x}^{(k+1)} &= -\mathbf{f}^{(k)}, \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k+1)}, \\ \Delta \mathbf{f}^{(k+1)} &= \mathbf{f}(\mathbf{x}^{(k+1)}) - \mathbf{f}(\mathbf{x}^{(k)}), \\ \mathbf{B}^{(k+1)} &= \mathbf{B}^{(k)} + \frac{1}{\|\Delta \mathbf{x}^{(k+1)}\|_2^2} (\Delta \mathbf{f}^{(k+1)} - \mathbf{B}^{(k)} \Delta \mathbf{x}^{(k+1)}) (\Delta \mathbf{x}^{(k+1)})^T\end{aligned}$$

Этот метод обладает сверхлинейной сходимостью.

§ 7.5. О некоторых подходах к решению задач локализации и отыскания решений систем нелинейных уравнений

Подчеркнем еще раз, что одной из наиболее трудных проблем, возникающих при решении систем нелинейных уравнений, является задача локализации решения. Изложим некоторые подходы к ее решению, широко применяемые в практике вычислений. Указанные методы можно применять и для вычисления решения с заданной точностью ε .

1. Использование методов минимизации. Иногда бывает полезно свести задачу отыскания решения системы нелинейных уравнений к задаче отыскания минимума функции многих переменных. В простейшем варианте это сведение выглядит следующим образом. Введем функцию $\Phi(x) = \sum_{i=1}^m (f_i(x))^2$. Она неотрицательна и достигает своего минимума (равного нулю) тогда и только тогда, когда $f_i(x) = 0$ для всех $i = 1, 2, \dots, m$, т.е. когда x является решением системы (7.1).

Применив для отыскания минимума функции Φ один из итерационных методов минимизации — методов спуска (см. гл. 10), можно найти вполне удовлетворительное приближение к решению \bar{x} , которое затем имеет смысл использовать как начальное приближение $x^{(0)}$ в одном из итерационных методов решения нелинейных систем.

Выгода от указанного сведения исходной задачи к задаче минимизации состоит в том, что, как правило, методы спуска имеют более широкую область сходимости. Использование методов спуска можно рассматривать и как один из способов локализации решений системы (7.1). Применение на заключительном этапе методов, специально ориентированных на решение нелинейных систем, вызвано тем, что вблизи искомого решения методы спуска сходятся медленнее.

Следует отметить, что функция $\Phi(x)$ может иметь и ненулевые локальные минимумы, и в зависимости от выбора начального приближения методы спуска могут приводить к точкам локального минимума, не являющимся решениями системы (7.1).

Пример 7.4. Решения системы (7.4) являются точками глобального минимума функции

$$\Phi(x_1, x_2) = (x_1^3 + x_2^3 - 8x_1x_2)^2 + (x_1 \ln x_2 - x_2 \ln x_1)^2. \blacksquare$$

2. Метод продолжения по параметру. Сначала заметим, что довольно часто на практике встречаются ситуации, когда система нелинейных уравнений естественным образом зависит от некоторого параметра t , т.е. имеет вид

$$f(x, t) = 0, \quad (7.24)$$

а решение ее следует найти при некотором фиксированном значении параметра, например при $t = 1$. В частности, этим свойством обладает система из примера 7.1. Предположим, что при каждом $t \in [0, 1]$ система (7.24) имеет решение $\bar{x} = \bar{x}(t)$, непрерывно зависящее от

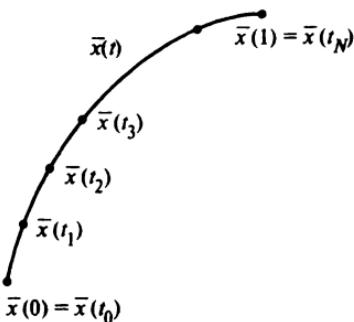


Рис. 7.5

параметра t , причем при $t = 0$ решение системы $f(x, t) = 0$ известно либо легко вычисляется.

Таким образом, семейство решений $\bar{x}(t)$ описывает в пространстве \mathbb{R}^m траекторию, начальная точка $\bar{x}(0)$ которой известна, а конечная $\bar{x}(1)$ подлежит определению (рис. 7.5).

И в тех случаях, когда система не зависит от параметра, можно ввести параметр t так, чтобы были выполнены указанные выше условия. Например, если x^* — известное приближение к решению системы $f(x) = 0$, то можно рассмотреть систему вида (7.24), положив

$$f(x, t) = f(x) - (1 - t)f(x^*).$$

Введем на отрезке $[0, 1]$ набор точек $0 = t_0 < t_1 < \dots < t_N = 1$. Используя тот или иной итерационный метод, решим последовательно для $k = 0, 1, 2, \dots, N$ системы $f(x, t_k) = 0$. При этом за начальное приближение к $\bar{x}(t_k)$ будем принимать решение $\bar{x}(t_{k-1})$. Если разность $t_k - t_{k-1}$ достаточно мала, то можно ожидать, что $\bar{x}(t_{k-1})$ будет достаточно хорошим начальным приближением к $\bar{x}(t_k)$, обеспечивающим сходимость используемого итерационного метода.

Замечание. Довольно часто на практике проводится исследование зависимости определенных характеристик объекта от некоторого параметра t . Для таких задач метод продолжения естествен. Более того, точки t_0, t_1, \dots, t_N можно выбрать, использовав дополнительные соображения, причем решение $\bar{x}(t_k)$ представляет интерес не только для $k = N$, но и для всех $k = 0, 1, \dots, N$.

3. Метод дифференцирования по параметру. Предположим, что решение $\bar{x}(t)$ системы (7.24) является гладкой функцией параметра t . Дифференцируя тождество $f(\bar{x}(t), t) = 0$ по t , получаем

$$f'_x(\bar{x}(t), t)\bar{x}'(t) + f'_t(\bar{x}(t), t) = 0.$$

Здесь $f'_x(\bar{x}, t)$ — матрица с элементами $\partial f_i / \partial x_j(x, t)$ ($i, j = 1, 2, \dots, m$), а $f'_t(\bar{x}, t)$ — вектор-столбец с элементами $\partial f_i / \partial t(x, t)$ ($i = 1, 2, \dots, m$).

Таким образом, если матрица f'_x невырождена ($\det f'_x \neq 0$), то $\bar{x}(t)$ является решением задачи Коши для системы обыкновенных дифференциальных уравнений

$$\dot{x}'(t) = -[f'_x(x(t), t)]^{-1}f'_t(x(t), t), \quad x(0) = \bar{x}(0). \quad (7.25)$$

Интересующее нас значение решения $\bar{x}(1)$ можно теперь найти приближенно, применяя численные методы решения задачи Коши (см. гл. 14).

Например, метод Эйлера приводит к процессу

$$x^{(k+1)} = x^{(k)} - (t_{k+1} - t_k)[f'_x(x^{(k)}, t_k)]^{-1}f'_t(x^{(k)}, t_k), \quad k = 0, 1, \dots, N-1,$$

где $x^{(0)} = \bar{x}(0)$, $x^{(k)}$ — приближение к $\bar{x}(t_k)$ (рис. 7.6).

Конечно, метод Эйлера приведен здесь только для удобства иллюстрации, а в реальной ситуации используется один из методов более высокого порядка точности.

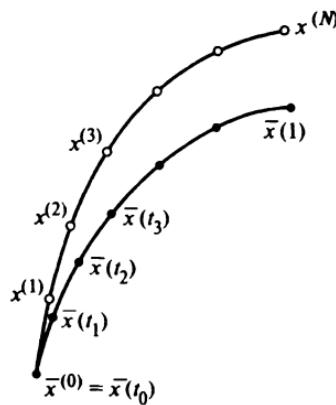


Рис. 7.6

Полученное указанным способом значение $x^{(N)}$ можно использовать и как хорошее начальное приближение к $\bar{x}(1)$ в одном из итерационных методов решения системы (7.24).

Замечание 1. Иногда метод дифференцирования по параметру называют *методом Давиденко*¹.

Замечание 2. Методы продолжения и дифференцирования по параметру нередко позволяют успешно преодолевать непростую проблему локализации. Однако следует отметить, что эти методы далеко не всегда оказываются эффективными и их практическое применение требует определенной осторожности.

§ 7.6. Дополнительные замечания

1. Существенно более подробная и богатая информация о системах нелинейных уравнений и методах их решения содержится в книгах [74, 40]. Рекомендуем также обратиться к учебнику [9].

2. В последнее время значительный интерес проявляется к так называемым квазиньютоновским методам, которые получаются при специальных аппроксимациях матрицы Якоби. Часто такая аппроксимация сочетается с использованием гибридных алгоритмов подобно тому, как это делается для решения одного нелинейного уравнения (см. § 4.8). Полезные обсуждения таких методов можно найти в [40].

3. Иногда подходящим для решения задачи методом оказывается метод установления см. [9]. Чаще это бывает тогда, когда решение системы (7.1) описывает устойчивое стационарное состояние некоторой физической системы.

4. Рекомендовать тот или иной метод как наиболее подходящий для решения системы нелинейных уравнений невозможно без тщательного анализа конкретной задачи.

¹ Дмитрий Федорович Давиденко (1923—1995) — российский математик.

Глава 8

МЕТОДЫ РЕШЕНИЯ ПРОБЛЕМЫ СОБСТВЕННЫХ ЗНАЧЕНИЙ

Вычисление собственных значений (чисел) и собственных векторов — одна из тех сложных вычислительных задач, с которой часто приходится сталкиваться исследователю, занимающемуся конструированием или анализом больших технических систем. В электрических и механических системах собственные числа отвечают собственным частотам колебаний, а собственные векторы характеризуют соответствующие формы (моды) колебаний. Знание собственных чисел позволяет анализировать многие процессы, исследовать и управлять ими. Оценка критических нагрузок при расчете строительных конструкций также основана на информации о собственных значениях и собственных векторах матриц.

Собственные числа и собственные векторы являются важнейшими характеристиками, отражающими существенные стороны линейных моделей. Поэтому, конечно, дальнейшее расширение процесса математического моделирования приведет к тому, что владение методами решения проблемы собственных значений станет неотъемлемым элементом высшего образования.

§ 8.1. Постановка задачи.

Некоторые вспомогательные сведения

1. Постановка задачи. В данной главе мы ограничимся рассмотрением методов решения проблемы собственных значений только для квадратных матриц A порядка m с вещественными элементами a_{ij} , ($i, j = 1, 2, \dots, m$). Будем всюду под $\|x\|$ понимать норму $\|x\|_2$ и под (x, y) — скалярное произведение векторов x, y (см. § 5.2).

Напомним, что число λ называется *собственным значением (собственным числом)* матрицы A , если существует ненулевой вектор x , удовлетворяющий системе уравнений

$$Ax = \lambda x \tag{8.1}$$

и называемый *собственным вектором* матрицы A , отвечающим собственному значению λ . (Собственное значение λ и собственный вектор x — вообще говоря, — комплексные.)

Запишем систему (8.1) в виде

$$(\mathbf{A} - \lambda \mathbf{E})\mathbf{x} = 0.$$

Эта однородная система имеет ненулевое решение \mathbf{x} тогда и только тогда, когда определитель матрицы системы равен нулю, т.е.

$$\det(\mathbf{A} - \lambda \mathbf{E}) = 0. \quad (8.2)$$

Раскрытие этого уравнения приводит к так называемому *характеристическому (или вековому) уравнению*

$$\lambda^m + p_1 \lambda^{m-1} + p_2 \lambda^{m-2} + \dots + p_{m-1} \lambda + p_m = 0, \quad (8.3)$$

представляющему собой алгебраическое уравнение степени m .

Известно, что характеристическое уравнение имеет в области комплексных чисел ровно m корней $\lambda_1, \lambda_2, \dots, \lambda_m$ (с учетом их кратности). Таким образом, каждая квадратная матрица \mathbf{A} порядка m обладает набором из m собственных значений $\lambda_1, \lambda_2, \dots, \lambda_m$.

Если матрица \mathbf{A} симметрична, то все ее собственные значения являются вещественными числами. Для несимметричных матриц возможно наличие комплексных собственных значений вида $\lambda = \alpha + i\beta$ с ненулевой мнимой частью. В этом случае собственным значением обязательно является и комплексно-сопряженное число $\bar{\lambda} = \alpha - i\beta$.

В ряде задач механики, физики, химии, техники, биологии требуется получение всех собственных значений некоторых матриц, а иногда и всех собственных векторов. В такой постановке задачу называют *полной проблемой собственных значений*.

Довольно часто определению подлежат не все собственные значения и собственные векторы, а лишь небольшая их часть. Например, существенный интерес во многих приложениях представляют максимальное или минимальное по модулю собственное значение или же собственное значение, наиболее близко расположенное к заданному значению. Такие задачи являются примерами *частичных проблем собственных значений*.

Может показаться, что достаточно ограничиться только рассмотрением методов решения полной проблемы собственных значений, так как все остальные проблемы являются ее частными случаями. Однако такой подход неоправдан, поскольку ориентирует на работу по получению значительного объема заведомо ненужной информации и требует существенно большего объема вычислений, чем это необходимо в действительности. Поэтому для решения различных частичных проблем собственных значений разработан ряд специальных методов.

Пример 8.1. Найдем собственные числа матрицы

$$A = \begin{bmatrix} 2 & -9 & 5 \\ 1.2 & -5.3999 & 6 \\ 1 & -1 & -7.5 \end{bmatrix}. \quad (8.4)$$

Запишем характеристический многочлен

$$\begin{aligned} P_3(\lambda) &= \det(A - \lambda E) = \det \begin{bmatrix} 2 - \lambda & -9 & 5 \\ 1.2 & -5.3999 - \lambda & 6 \\ 1 & -1 & -7.5 - \lambda \end{bmatrix} = \\ &= -\lambda^3 - 10.8999\lambda^2 - 26.49945\lambda - 21.002. \end{aligned}$$

Используя один из итерационных методов решения нелинейных уравнений (например, метод Ньютона), определяем один из корней уравнения $P_3(\lambda) = 0$, а именно $\lambda_1 \approx -7.87279$. Разделив $P_3(\lambda)$ на $\lambda - \lambda_1$, имеем

$$\frac{P_3(\lambda)}{\lambda + 7.87279} \approx P_2(\lambda) = \lambda^2 + 3.02711\lambda + 2.66765.$$

Решая квадратное уравнение $P_2(\lambda) = 0$, находим корни $\lambda_{2,3} \approx -1.51356 \pm 0.613841 i$.

Таким образом, матрица A имеет одно вещественное собственное значение $\lambda_1 \approx -7.87279$ и два комплексно-сопряженных собственных значения $\lambda_{2,3} \approx -1.51356 \pm 0.613841 i$. ▲

Численные методы решения проблемы собственных значений, использовавшиеся до конца 40-х годов XX в. сводились в конечном счете к решению характеристического уравнения (8.3). Этой классической схеме следовали и мы в примере 8.1. При реализации такого подхода основные усилия были направлены на разработку эффективных методов быстрого вычисления коэффициентов характеристического уравнения. Методы такого класса получили названия *прямых*; к ним относятся пользовавшиеся популярностью методы Крылова¹, Данилевского, Леверье и др.

Однако указанный подход становится неудовлетворительным, если речь идет о вычислении собственных значений матриц, имеющих порядок m в несколько десятков (и тем более сотен), т.е. матриц довольно скромных по современным понятиям размеров.

¹ Алексей Николаевич Крылов (1863—1945) — русский математик, механик и кораблестроитель, является автором первого в мировой научной литературе курса по численным методам — изданных в 1911 г. «Лекций о приближенных вычислениях».

Одна из причин состоит в том, что хотя задачи (8.1) и (8.3) формально эквивалентны, они имеют разную обусловленность. Так как корни многочлена $P_m(x)$ высокой степени чрезвычайно чувствительны к погрешностям в коэффициентах, то на этапе вычисления коэффициентов характеристического уравнения может быть в значительной степени потеряна информация о собственных значениях матрицы.

С появлением компьютеров широкое распространение получили итерационные методы решения проблемы собственных значений, не использующие вычисление характеристического многочлена. К началу 60-х годов XX в. эти методы практически полностью вытеснили прямые методы из практики вычислений.

2. Преобразование подобия. Говорят, что матрицы A и B подобны, если существует невырожденная матрица P (*матрица подобия*) такая, что $B = P^{-1}AP$. Само преобразование матрицы A к виду $B = P^{-1}AP$ называют *преобразованием подобия*. Преобразование подобия матрицы возникает естественным образом как результат замены переменных (или перехода к новому базису) в пространстве m -мерных векторов. Действительно, пусть y — результат применения матрицы A к вектору x , т.е. $y = Ax$. Произведем замену переменных $x = Px'$, $y = Py'$. Тогда равенство $y = Ax$ примет вид $y' = P^{-1}APx'$. Это означает, что в новых переменных то же самое преобразование осуществляется уже не матрицей A , а матрицей $P^{-1}AP$, подобной A .

Важно то, что и полученная в результате преобразования подобия матрица имеет тот же набор собственных чисел. В самом деле, рассмотрим характеристический многочлен для матрицы $P^{-1}AP$ и воспользуемся тем, что определитель произведения квадратных матриц равен произведению соответствующих определителей:

$$\det(P^{-1}AP - \lambda E) = \det(P^{-1}(A - \lambda E)P) = \det P^{-1} \det(A - \lambda E) \det P = \det(A - \lambda E).$$

Таким образом, характеристические многочлены, а следовательно, и собственные числа матриц A и $P^{-1}AP$ совпадают. Соответствующие собственные векторы x и x' не совпадают, но, как нетрудно установить, они связаны равенством $x = Px'$.

Если бы матрицу A с помощью преобразования подобия или последовательности таких преобразований удалось привести к верхнему треугольному виду, то проблему вычисления собственных значений

можно было бы считать решенной. Дело в том, что у верхней треугольной матрицы

$$\mathbf{B} = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1m} \\ 0 & b_{22} & \dots & b_{2m} \\ \dots & & & \\ 0 & 0 & \dots & b_{mm} \end{bmatrix} \quad (8.5)$$

собственными числами являются элементы главной диагонали b_{ii} . В этом нетрудно убедиться, если записать характеристический многочлен: $\det(\mathbf{B} - \lambda E) = (b_{11} - \lambda)(b_{22} - \lambda) \dots (b_{mm} - \lambda)$.

Оказывается, что преобразованием подобия матрицу A можно привести к еще более простому виду, чем (8.5). Справедлива (см. [27]) следующая теорема.

Теорема 8.1. Любой квадратной матрицы A с помощью преобразования подобия можно привести к следующему виду:

$$\mathbf{P}^{-1} \mathbf{A} \mathbf{P} = \mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & \sigma_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & \lambda_2 & \sigma_2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \lambda_3 & \sigma_3 & \dots & 0 & 0 \\ \dots & & & & & & \\ 0 & 0 & 0 & 0 & \dots & \lambda_{m-1} & \sigma_{m-1} \\ 0 & 0 & 0 & 0 & \dots & 0 & \lambda_m \end{bmatrix}. \quad (8.6)$$

Здесь $\lambda_1, \lambda_2, \dots, \lambda_m$ — собственные числа матрицы A . Числа σ_i принимают одно из двух значений 0 или 1, причем если $\sigma_i = 1$, то обязательно $\lambda_i = \lambda_{i+1}$. ■

Матрица (8.6) называется *жордановой формой* матрицы A .

3. Матрицы простой структуры. В этой главе особое внимание будет уделено *матрицам простой структуры*, т.е. матрицам, которые с помощью преобразования подобия можно привести к диагональному виду:

$$\mathbf{P}^{-1} \mathbf{A} \mathbf{P} = \mathbf{D} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \dots & & & \\ 0 & 0 & \dots & \lambda_m \end{bmatrix}. \quad (8.7)$$

Запишем равенство (8.7) в виде $AP = PD$ и заметим, что оно верно тогда и только тогда, когда каждый j -й столбец матрицы P является собственным вектором матрицы A , отвечающим собственному значению λ_j . Таким образом, верна следующая теорема.

Теорема 8.2. *Матрица A является матрицей простой структуры тогда и только тогда, когда она имеет m линейно независимых собственных векторов e_1, e_2, \dots, e_m , отвечающих собственным значениям $\lambda_1, \lambda_2, \dots, \lambda_m$ соответственно.* ■

Указанные в теореме 8.2 собственные векторы e_1, e_2, \dots, e_m , образуют базис в пространстве m -мерных векторов. Поэтому каждый m -мерный вектор x может быть однозначно представлен в виде линейной комбинации этих векторов:

$$x = c_1 e_1 + c_2 e_2 + \dots + c_m e_m. \quad (8.8)$$

Какие же матрицы заведомо могут быть приведены к диагональному виду? Следующие два предложения частично отвечают на этот вопрос.

Теорема 8.3. *Если все собственные значения матрицы A различны, то она является матрицей простой структуры.* ■

Теорема 8.4. *Если A — вещественная симметричная матрица, то, она подобна диагональной матрице, причем матрица подобия P может быть выбрана ортогональной (т.е. удовлетворяющей условию $P^{-1} = P^T$).* ■

4. Локализация собственных значений. С помощью так называемых *теорем локализации* иногда удается получить грубые оценки расположения собственных чисел. Изложим самый известный из этих результатов — теорему Гершгорина¹.

Пусть $r_i = \sum_{\substack{j=1 \\ j \neq i}}^m |a_{ij}|$ — сумма модулей внедиагональных элементов i -й строки матрицы A . Обозначим через S_i замкнутый круг радиуса r_i на комплексной плоскости с центром в точке a_{ii} , т.е. $S_i = \{z \in \mathbb{C}: |z - a_{ii}| \leq r_i\}$. Будем называть круги S_i *кругами Гершгорина*. Имеет место следующее утверждение.

¹ Семен Аронович Гершгорин (1901—1933) — российский математик. Результат, о котором идет речь, был опубликован им в 1931 г.

Теорема 8.5. (теорема Гершгорина). Все собственные значения матрицы A лежат в объединении кругов S_1, S_2, \dots, S_m .

Доказательство. Возьмем произвольное собственное значение λ матрицы A и соответствующий собственный вектор x . Пусть x_i — максимальная по модулю координата вектора x . Запишем i -е уравнение системы (8.1) в следующем виде:

$$(a_{ii} - \lambda)x_i = -\sum_{j \neq i} a_{ij}x_j.$$

Из этого равенства с учетом оценки $|x_j/x_i| \leq 1$ следует неравенство

$$|a_{ii} - \lambda| \leq \sum_{j \neq i} |a_{ij}| \left| \frac{x_j}{x_i} \right| \leq r_i.$$

Таким образом, $\lambda \in S_i$. ■

Пример 8.2. Для матрицы

$$A = \begin{bmatrix} -2 & 0.5 & 0.5 \\ -0.5 & -3.5 & 1.5 \\ 0.8 & -0.5 & 0.5 \end{bmatrix}$$

круги Гершгорина изображены на рис. 8.1. Здесь $r_1 = 1$, $r_2 = 2$, $r_3 = 1.3$. ▲

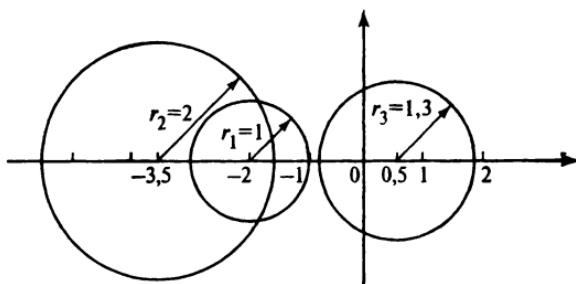


Рис. 8.1

Следующий результат является полезным дополнением к теореме Гершгорина.

Теорема 8.6. Если k кругов Гершгорина образуют замкнутую область \bar{G} , изолированную от других кругов, то в области \bar{G} находится ровно k собственных значений матрицы A (с учетом их кратности).

Следствие. Если какой-либо круг Гершгорина изолирован, то он содержит ровно одно собственное значение матрицы A . ■

Пример 8.3. Для матрицы A из примера 8.2 в объединении кругов S_1 и S_2 находится ровно два собственных значения λ_1 и λ_2 , а круг S_3 содержит ровно одно собственное значение λ_3 . ▲

5. Отношение Рэлея. При вычислении собственных чисел и собственных векторов симметричных матриц важную роль играет функция

$$\rho(x) = \frac{(Ax, x)}{(x, x)}, \quad (8.9)$$

называемая *отношением Рэлея*¹. Следующее предложение частично объясняет значение этой величины.

Теорема 8.7. Пусть A — симметричная матрица. Тогда справедливы следующие утверждения:

1) минимальное и максимальное собственные значения матрицы A вычисляются по формулам

$$\lambda_{\min} = \min_{x \neq 0} \rho(x), \quad \lambda_{\max} = \max_{x \neq 0} \rho(x);$$

2) вектор x является стационарной точкой функции $\rho(x)$ (т.е. удовлетворяет условию $\operatorname{grad} \rho(x) = 0$) тогда и только тогда когда x — собственный вектор матрицы A . ■

При построении методов решения проблемы собственных значений существенно используется то обстоятельство, что если x — хорошее приближение к собственному вектору e_j , то $\rho(x)$ — хорошее приближение к собственному числу λ_j .

6. Обусловленность задачи вычисления собственных значений и собственных векторов. Пусть A_* — матрица с приближенно заданными элементами $a_{ij}^* \approx a_{ij}$. Обозначим через λ_j^* ($j = 1, 2, \dots, m$) собственные числа матрицы A_* . Рассмотрение вопроса о том, как влияет погрешность задания матрицы на погрешность искомых собственных значений, начнем с формулировки следующего известного результата.

Теорема 8.8. Пусть A и A_* — симметричные матрицы, а λ_j и λ_j^* — их собственные числа, упорядоченные по возрастанию. Тогда справедливы оценки погрешности

$$\max_{1 \leq j \leq m} |\lambda_j - \lambda_j^*| \leq \|A - A_*\|_2, \quad (8.10)$$

$$\left(\sum_{j=1}^m (\lambda_j - \lambda_j^*)^2 \right)^{1/2} \leq \|A - A_*\|_E. \quad (8.11)$$

¹Джон Уильям Рэлей(1842—1919) — английский физик и математик.

Теорема 8.8 означает, что задача вычисления собственных значений симметричных матриц очень хорошо обусловлена. Следовательно, в этом случае собственные числа надежно определяются заданием элементов матрицы. К сожалению, для несимметричных матриц дело обстоит совсем иначе. Хотя задача вычисления собственных значений и в этом случае является устойчивой, для многих несимметричных матриц собственные значения чрезвычайно чувствительны к погрешностям задания коэффициентов.

Пример 8.4. Приведем принадлежащий Дж.Х. Уилкинсону [102] пример матрицы, очень плохо обусловленной по отношению к проблеме собственных значений.

Собственными числами верхней треугольной матрицы 20-го порядка

$$\mathbf{A} = \begin{bmatrix} 20 & 20 & 0 & 0 & \dots & 0 & 0 \\ 0 & 19 & 20 & 0 & \dots & 0 & 0 \\ 0 & 0 & 18 & 20 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 2 & 20 \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 \end{bmatrix}$$

являются числа $\lambda_1 = 1$, $\lambda_2 = 2$, ..., $\lambda_{20} = 20$. Заметим, что для этой матрицы характеристический многочлен $P_{20}(\lambda) = (20 - \lambda)(19 - \lambda) \dots (1 - \lambda)$ совпадает с многочленом, рассмотренным в примере 3.8 в связи с плохой обусловленностью его корней. Добавим малое ε к элементу $a_{20,1} = 0$. В результате характеристическое уравнение примет вид

$$(20 - \lambda)(19 - \lambda) \dots (1 - \lambda) = 20^{19}\varepsilon.$$

При $\varepsilon = 10^{-10}$ собственные значения возмущенной матрицы таковы: $\lambda_1^* \approx 0.996$, $\lambda_2^* \approx 2.11$, $\lambda_3^* \approx 2.57$, $\lambda_{4,5}^* \approx 3.97 \pm 1.09i$, $\lambda_{6,7}^* \approx 5.89 \pm 1.95i$, $\lambda_{8,9}^* \approx 8.12 \pm 2.53i$, $\lambda_{10,11}^* \approx 10.5 \pm 2.73i$, $\lambda_{12,13}^* \approx 12.9 \pm 2.53i$, $\lambda_{14,15}^* \approx 15.1 \pm 1.95i$, $\lambda_{16,17}^* \approx 17.0 \pm 1.09i$, $\lambda_{18}^* \approx 18.4$, $\lambda_{19}^* \approx 18.9$, $\lambda_{20}^* \approx 20.0$. Как нетрудно видеть, большинство значений оказались полностью искаженными. ▲

Отметим, что число обусловленности $\text{cond}(\mathbf{A})$ не характеризует обусловленность матрицы \mathbf{A} по отношению к проблеме собственных значений. Оказывается, что такой характеристикой чувствительности собственных значений относительно погрешности задания матрицы для

матрицы простой структуры служит число обусловленности матрицы P , столбцы которой являются собственными векторами матрицы A . В подтверждение сказанного приведем следующий результат.

Теорема 8.9. Пусть $P^{-1}AP = D$, где D — диагональная матрица из собственных значений матрицы A . Тогда каждое собственное значение матрицы A_* удалено от некоторого собственного значения матрицы A не более чем на $d = \text{cond}_2(P) \|A - A_*\|_2$. ■

Пусть x — собственный вектор матрицы A , отвечающий собственному значению λ , а x^* — собственный вектор приближенно заданной матрицы A_* , отвечающий собственному значению λ^* . Прежде чем обсуждать обусловленность задачи вычисления собственного вектора x , заметим, что здесь вопрос о выборе меры близости векторов x^* и x не является тривиальным. Выбор в качестве такой меры величины $\|x - x^*\|$ неудачен. Дело в том, что собственные векторы x^* и x не определяются однозначно. Во всяком случае после умножения каждого из них на любые числа $\alpha_1 \neq 0$ и $\alpha_2 \neq 0$ полученные векторы $\alpha_1 x^*$ и $\alpha_2 x$ снова являются собственными. Поэтому имеет смысл стремиться не к тому, чтобы векторы x^* и x были близки по норме, а к тому, чтобы они были близки по направлению. Исходя из этого, примем в качестве меры близости x^* к x величину $|\sin \varphi|$, где φ — угол между векторами x^* и x , вычисляемый по формуле

$$\varphi = \arccos \left(\frac{(x^*, x)}{\|x^*\|_2 \|x\|_2} \right).$$

Задача вычисления собственных векторов симметричной матрицы хорошо обусловлена, если собственные значения хорошо отделены друг от друга. В подтверждение сказанного приведем следующий результат.

Теорема 8.10. Пусть A и A_* — симметричные матрицы. Тогда верна оценка

$$|\sin \varphi| \leq \frac{\|A - A_*\|_2}{\gamma}.$$

Здесь φ — угол между векторами x^* и x , а γ — расстояние от λ^* до ближайшего из несовпадающих с λ собственных значений матрицы A . ■

В случае, когда матрица A несимметрична, собственные векторы могут оказаться очень плохо обусловленными.

Замечание. Подавляющее большинство встречающихся на практике матриц являются матрицами простой структуры. Это обстоятельство, а также большая простота анализа методов вычислений и формулировок соответствующих результатов позволяют нам в основном остановиться на проблеме собственных значений для матриц простой структуры.

§ 8.2. Степенной метод

Пусть требуется вычислить максимальное по модулю собственное значение λ_1 матрицы A , причем известно, что

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_m|. \quad (8.12)$$

Замечание. Собственное значение λ_1 в данном случае должно быть вещественным, поскольку в противном случае собственным значением было бы также равное ему по модулю число $\bar{\lambda}_1$.

1. Степенной метод без сдвигов. Опишем простейший вариант степенного метода, применяемого для вычисления λ_1 . Возьмем произвольный начальный вектор $x^{(0)} \neq 0$ и построим последовательности векторов $\{x^{(k)}\}_{k=0}^{\infty}$ и приближений $\{\lambda_1^{(k)}\}_{k=1}^{\infty}$ к λ_1 , использовав формулы

$$x^{(k)} = Ax^{(k-1)}, \quad (8.13)$$

$$\lambda_1^{(k)} = \frac{(x^{(k)}, x^{(k-1)})}{(x^{(k-1)}, x^{(k-1)})}. \quad (8.14)$$

Замечание. Правая часть формулы (8.14) — это просто отношение Рэлея, вычисленное при $x = x^{(k-1)}$. В самом деле, $\lambda_1^{(k)} = (Ax^{(k-1)}, x^{(k-1)})/(x^{(k-1)}, x^{(k-1)}) = \rho(x^{(k-1)})$.

Справедлива следующая теорема.

Теорема 8.11. Пусть A — матрица простой структуры, для которой выполнено условие (8.12). Предположим, что в разложении

$$x^{(0)} = c_1 e_1 + c_2 e_2 + \dots + c_m e_m = \sum_{i=1}^m c_i e_i, \quad (8.15)$$

по нормированному базису из собственных векторов коэффициент c_1 не равен нулю. Тогда $\lambda_1^{(k)} \rightarrow \lambda_1$ при $k \rightarrow \infty$ и при всех $k \geq k_0(x^{(0)})$ верна оценка относительной погрешности

$$\delta(\lambda_1^{(k)}) = \frac{|\lambda_1^{(k)} - \lambda_1|}{|\lambda_1|} \leq C_0 \left| \frac{\lambda_2}{\lambda_1} \right|^{k-1}, \quad (8.16)$$

где $C_0 = 4 \sum_{i=2}^m |\alpha_i|$, а $\alpha_i = c_i/c_1$, $2 \leq i \leq m$.

Доказательство. Как нетрудно видеть, вектор $x^{(k)}$ удовлетворяет равенству $x^{(k)} = A^k x^{(0)}$, т.е. получается из $x^{(0)}$ умножением на k -ю степень матрицы A (отсюда и название метода). Так как $A e_i = \lambda_i e_i$, то $A^k e_i = \lambda_i^k e_i$ и из (8.15) следует

$$x^{(k)} = \lambda_1^k c_1 e_1 + \sum_{i=2}^m \lambda_i^k c_i e_i. \quad (8.17)$$

Положим $\mu_i = \lambda_i / \lambda_1$, $2 \leq i \leq m$ и

$$z^{(k)} = \frac{x^{(k)}}{\lambda_1^{(k)} c_1} = e_1 + \sum_{i=2}^m \mu_i^k \alpha_i e_i \quad (8.18)$$

Заметим, что

$$z^{(k)} - z^{(k-1)} = \sum_{i=2}^m (\mu_i - 1) \mu_i^{k-1} \alpha_i e_i. \quad (8.19)$$

Так как $|\mu_i| \leq |\mu_2| < 1$ для всех $2 \leq i \leq m$, а базис нормированный (т.е. $\|e_i\|_2 = 1$ для всех i), то из (8.18), (8.19) следует, что

$$\|z^{(k)}\|_2 \geq \|e_1\|_2 - \sum_{i=2}^m |\mu_i|^k |\alpha_i| \|e_i\|_2 \geq 1 - |\mu_2|^k \sum_{i=2}^m |\alpha_i| = 1 - \frac{1}{4} C_0 |\mu_2|^k, \quad (8.20)$$

$$\|z^{(k)} - z^{(k-1)}\|_2 \leq 2 |\mu_2|^{k-1} \sum_{i=2}^m |\alpha_i| \|e_i\|_2 = 2 |\mu_2|^{k-1} \sum_{i=2}^m |\alpha_i| = \frac{1}{2} C_0 |\mu_2|^{k-1}.$$

Выберем $k_0 = k_0(x^{(0)})$ таким, чтобы $|\mu_2|^{k_0-1} C_0 \leq 2$. Тогда для всех $k \geq k_0$ в силу неравенства (8.20) будет справедливо неравенство $\|\mathbf{z}^{(k-1)}\|_2 \geq 1/2$. Осталось заметить, что

$$\begin{aligned}\delta(\lambda_1^{(k)}) &= \left| \frac{\lambda_1^{(k)}}{\lambda_1} - 1 \right| = \left| \frac{(x^{(k)}, x^{(k-1)})}{\lambda_1(x^{(k-1)}, x^{(k-1)})} - 1 \right| = \left| \frac{(z^{(k)}, z^{(k-1)})}{(z^{(k-1)}, z^{(k-1)})} - 1 \right| = \\ &= \frac{|(z^{(k)} - z^{(k-1)}, z^{(k-1)})|}{\|z^{(k-1)}\|_2^2} \leq \frac{\|z^{(k)} - z^{(k-1)}\|_2}{\|z^{(k-1)}\|_2} \leq C_0 |\mu_2|^{k-1}.\end{aligned}$$

(При выводе этой оценки мы воспользовались неравенством Коши—Буняковского $|(x, y)| \leq \|x\|_2 \|y\|_2$ для $x = z^{(k)} - z^{(k-1)}$ и $y = z^{(k-1)}$.) ■

Замечание 1. Если выполнены условия теоремы 18.11 и матрица A симметрична, то оценка (8.16) верна уже при всех $k \geq 1$

с постоянной $C_0 = \sqrt{\sum_{i=2}^m |\alpha_i|^2}$.

Замечание 2. При выполнении условия (8.12) оценка (8.16) (с другой постоянной C_0) верна и для произвольной матрицы, а не только для матрицы простой структуры.

Замечание 3. Если $|\lambda_1| > 1$, то $\|x^{(k)}\|_2 \rightarrow \infty$ при $k \rightarrow \infty$ (см. формулу (8.17) и при вычислении на компьютере возможно переполнение). Если же $|\lambda_1| < 1$, то $\|x^{(k)}\|_2 \rightarrow 0$ и возможно исчезновение порядка. Для предупреждения этих ситуаций обычно вектор $x^{(k)}$ нормируют так, чтобы $\|x^{(k)}\|_2 = 1$ для всех $k \geq 1$. Одна из используемых модификаций такова:

$$y^{(k)} = Ax^{(k-1)}, \quad \lambda^{(k)} = (y^{(k)}, x^{(k-1)}), \quad x^{(k)} = \frac{y^{(k)}}{\|y^{(k)}\|_2}. \quad (8.21)$$

Предполагается, что $\|x^{(0)}\|_2 = 1$.

Теорема 8.11 для метода (8.21) остается справедливой. Более того, последовательность $x^{(k)}$ сходится к собственному вектору по направлению, т.е. $\sin \varphi^{(k)} \rightarrow 0$ при $k \rightarrow \infty$, где $\varphi^{(k)}$ — угол между векторами $x^{(k)}$ и e_1 .

Замечание 4. Крайне маловероятно, чтобы в разложении (8.15) вектора $x^{(0)}$ по собственным векторам коэффициент c_1 при e_1 оказался равным нулю. Теоретически в этой исключительной ситуации метод не должен давать сходящуюся к λ_1 последовательность. Однако при вычислении на компьютере по формулам (8.12) через несколько итераций вследствие погрешностей округления почти наверняка появится ненулевой коэффициент при e_1 в разложении вектора $x^{(k)}$. Итерационный процесс снова окажется сходящимся (правда, потребуется выполнить большее число итераций).

2. Апостериорная оценка погрешности. В общем случае для решения проблемы собственных значений не существует эффективных апостериорных оценок погрешности, в особенности для собственных векторов. Более простой эта проблема выглядит для симметричных матриц.

Для оценки погрешности приближения λ^* к собственному значению λ симметричной матрицы полезно использовать следующую теорему.

Теорема 8.12. Пусть λ^* — произвольное число, а x^* — произвольный ненулевой вектор. Тогда для любой симметричной матрицы A существует собственное значение этой матрицы λ такое, что справедлива оценка

$$|\lambda - \lambda^*| \leq \frac{\|Ax^* - \lambda^*x^*\|_2}{\|x^*\|_2}. \blacksquare \quad (8.22)$$

Если x^* — приближенно вычисленный собственный вектор, то неравенство (8.22) позволяет дать простую апостериорную оценку погрешности вычисленного собственного значения λ^* . В частности, для степенного метода (8.13), (8.14) из (8.22) следует такая апостериорная оценка:

$$|\lambda_1 - \lambda_1^{(k)}| \leq \frac{\|x^{(k)} - \lambda_1^{(k)}x^{(k-1)}\|_2}{\|x^{(k-1)}\|_2}$$

(напомним, что здесь предполагается симметричность матрицы A).

Пусть x^* — ненулевой вектор, рассматриваемый как приближение к собственному вектору x , который отвечает собственному значению λ . Далее, пусть $\rho^* = \rho(x^*)$ — отношение Рэлея, $r = Ax^* - \rho^*x^*$ — вектор невязки, $\sigma = \frac{\|r\|_2}{\|x^*\|_2}$. Если собственное значение λ хорошо отделено от остальных собственных значений матрицы A , то для оценки погрешности может быть использован следующий результат.

Теорема 8.13. Пусть λ — ближайшее к ρ^* собственное значение матрицы A , x — соответствующий собственный вектор и ϕ — угол между векторами x^* и x . Тогда справедливы оценки

$$|\sin \phi| \leq \frac{\sigma}{\gamma}, \quad |\lambda - \rho^*| \leq \frac{\sigma^2}{\gamma}.$$

Здесь $\gamma = \min_{\lambda_i \neq \lambda} |\lambda_i - \rho^*|$ — расстояние от ρ^* до ближайшего из отличных от λ собственных значений матрицы A . ■

Приведем еще одно важное свойство симметричных матриц.

Теорема 8.14. Пусть A — симметричная матрица, x^* — произвольный ненулевой вектор, $\rho^* = \rho(x^*)$. Тогда ρ^* и x^* являются точными собственным числом и собственным вектором некоторой симметричной матрицы A_* , для которой

$$\|A - A_*\|_2 = \sigma. \blacksquare \quad (8.23)$$

Оценка (8.23) может оказаться очень полезной для обратного анализа ошибок (о существе такого подхода к анализу ошибок мы говорили в § 3.6). Например, очень часто матрица A , для которой на компьютере производится вычисление собственных значений, является лишь приближенным представлением «истинной» матрицы \tilde{A} , собственными значениями и векторами которой в действительности и интересуется исследователь. Пусть известен порядок погрешности $\varepsilon = \|A - \tilde{A}\|_2$ приближенно заданной матрицы и для найденной

пары $x^*, \rho^* = \rho(x^*)$ величина $\sigma = \frac{\|Ax^* - \rho^*x^*\|_2}{\|x^*\|_2}$ оказывается сравнимой

с величиной ε или даже существенно меньше ее. В такой ситуации ρ^* и x^* оказываются собственными значением и вектором матрицы A_* , которая отличается от «истинной» матрицы \tilde{A} почти так же, как отличается от нее матрица A . Вряд ли тогда следует стремиться к получению «более точных» решений.

Пример 8.5. Используя степенной метод, найдем для матрицы (8.4) максимальное по модулю собственное значение λ_1 , а также отвечающий ему собственный вектор.

Возьмем $x^{(0)} = (1, 0, 0)^T$ и будем вести вычисления по формулам (8.21).

Первая итерация. Вычисляем

$$y_1^{(1)} = 2x_1^{(0)} - 9x_2^{(0)} + 5x_3^{(0)} = 2,$$

$$y_2^{(1)} = 1.2x_1^{(0)} - 5.3999x_2^{(0)} + 6x_3^{(0)} = 1.2,$$

$$y_3^{(1)} = x_1^{(0)} - x_2^{(0)} + 7.5x_3^{(0)} = 1.$$

Тогда

$$\lambda_1^{(1)} = (y^{(1)}, x^{(0)}) = y_1^{(1)}x_1^{(0)} + y_2^{(1)}x_2^{(0)} + y_3^{(1)}x_3^{(0)} = 2 \cdot 1 + 1.2 \cdot 0 + 1 \cdot 0 = 2.$$

Далее

$$\|y^{(1)}\|_2 = \sqrt{2^2 + 1.2^2 + 1^2} \approx 2.53772,$$

$$x^{(1)} = \frac{y^{(1)}}{\|y^{(1)}\|_2} \approx (0.788110, 0.472866, 0.394055)^T.$$

Результаты десяти первых итераций с шестью знаками мантиссы приведены в табл. 8.1.

Таблица 8.1

Номер итерации k	$\lambda_1^{(k)}$	$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$
0	—	1.000000	0.000000	0.000000
1	2.00000	0.788110	0.472866	0.394055
2	-1.24159	-0.250056	0.266743	-0.930763
3	-6.08121	-0.611733	-0.593143	0.523415
4	-9.50971	0.700557	0.583736	-0.410455
5	-8.23499	-0.716752	-0.579514	0.387853
6	-7.94591	0.720030	0.578443	-0.383353
7	-7.88567	-0.720635	-0.578216	0.382446
8	-7.87462	0.720729	0.578174	-0.382446
9	-7.87294	-0.720739	-0.578167	0.382437
10	-7.87277	-0.720739	0.578167	-0.382437

Хотя мы не имеем в данном случае обоснованного критерия окончания, по-видимому, при $k = 10$ следует прекратить вычисления и, округлив результаты, положить $\lambda_1 \approx -7.873$, $e_1 \approx (-0.7207, 0.5782, -0.3824)^T$. ▲

Важными достоинствами степенного метода являются его простота, возможность эффективного использования разреженности матрицы и отсутствие необходимости преобразования матрицы A .

Недостаток метода в применении к многим прикладным задачам — довольно медленная сходимость. Часто в приложениях значение $|\lambda_2|$

оказывается близким к $|\lambda_1|$. Так как скорость сходимости степенного метода определяется величиной $\left| \frac{\lambda_2}{\lambda_1} \right| \approx 1$, то в этом случае сходимость будет очень медленной.

3. Степенной метод со сдвигами. Существует несколько способов преодоления указанной трудности. Один из них заключается в применении степенного метода не к матрице A , а к матрице $\tilde{A} = A - \sigma E$, где σ — некоторое число. Собственными значениями матрицы \tilde{A} являются числа $\lambda_i - \sigma$, получаемые сдвигом собственных значений λ_i на число σ . Если число $\lambda_1 - \sigma$ по-прежнему остается максимальным по модулю, то следует попытаться подобрать σ так, чтобы сделать величину $\max_{2 \leq i \leq n} \left| \frac{\lambda_i - \sigma}{\lambda_1 - \sigma} \right|$ минимальной. Если, например все собственные значения положительны, то такой минимум достигается при $\sigma = (\lambda_2 + \lambda_m)/2$.

После того как приближение λ_1^* к первому собственному значению вычислено, степенной метод можно использовать для вычисления очередного собственного значения. Один из приемов такого использования состоит в сдвиге собственных значений на $\sigma = \lambda_1^*$. В этом случае число $\lambda_1 - \lambda_1^*$ станет минимальным по модулю, а максимальным по модулю окажется сдвиг другого собственного значения.

Существует несколько способов избавления от уже вычисленных собственных чисел и соответствующих собственных векторов с целью избежать их повторного вычисления. Эти способы принято называть *исчерпыванием*. Более подробную информацию о них можно найти, например в [23, 75].

§ 8.3. Метод обратных итераций

1. Вычисление собственных векторов методом обратных итераций. Многие из методов решения проблем собственных значений лучше приспособлены для вычисления собственных значений, чем собственных векторов. Поэтому целесообразно рассмотреть задачу вычисления собственного вектора e_j при условии, что уже найдено достаточно точное приближение λ_j^* к собственному значению λ_j .

Если исходить непосредственно из определения собственного вектора, то e_j следует искать как нетривиальное решение однородной системы уравнений

$$(A - \lambda_j E)x = 0 \quad (8.24)$$

с вырожденной матрицей $A - \lambda_j E$. Однако λ_j известно лишь приближенно и в действительности при таком подходе вместо системы (8.24) придется решать систему

$$(A - \lambda_j^* E)x = 0. \quad (8.25)$$

Так как матрица $A - \lambda_j^* E$ заведомо невырождена, то решением системы (8.25) является только $x = 0$. Следовательно, непосредственное численное решение системы (8.25) не дает возможность вычислить собственный вектор.

Одним из эффективных методов вычисления собственных векторов является *метод обратных итераций*. В этом методе приближение к собственному вектору определяют последовательным решением систем уравнений

$$(A - \lambda_j^* E)y^{(k+1)} = x^{(k)} \quad (8.26)$$

с последующей нормировкой решения:

$$x^{(k+1)} = \frac{y^{(k+1)}}{\|y^{(k+1)}\|_2}. \quad (8.27)$$

В качестве начального приближения берут ненулевой вектор $x^{(0)}$ с произвольно выбираемыми или даже случайными компонентами. Часто удовлетворительным является выбор $x^{(0)} = (1, 1, \dots, 1)^T$.

Чтобы понять механизм действия метода, рассмотрим случай, когда A — матрица простой структуры, а λ_j — простое собственное значение. Представим векторы $x^{(0)}$ и $y^{(1)}$ в виде линейных комбинаций собственных векторов e_1, e_2, \dots, e_m :

$$x^{(0)} = \sum_{i=1}^m c_i e_i, \quad y^{(1)} = \sum_{i=1}^m \alpha_i e_i.$$

Так как $(A - \lambda_j^* E)y^{(1)} = \sum_{i=1}^m \alpha_i (\lambda_i - \lambda_j^*) e_i$, то систему уравнений (8.26)

при $k = 0$ можно записать в виде

$$\sum_{i=1}^m \alpha_i (\lambda_i - \lambda_j^*) e_i = \sum_{i=1}^m c_i e_i.$$

Приравнивая коэффициенты при e_i , получаем $\alpha_i = \frac{c_i}{\lambda_i - \lambda_j^*}$. Следовательно,

$$\mathbf{y}^{(1)} = \sum_{i=1}^m \frac{c_i}{\lambda_i - \lambda_j^*} e_i = \frac{1}{\lambda_j - \lambda_j^*} \left[c_j e_j + \sum_{\substack{i=1 \\ i \neq j}}^m \frac{\lambda_j - \lambda_j^*}{\lambda_i - \lambda_j^*} c_i e_i \right]. \quad (8.28)$$

Если $|\lambda_j - \lambda_j^*| \ll |\lambda_i - \lambda_j^*|$ для всех $i \neq j$, то второе слагаемое в правой части формулы (8.28) мало по сравнению с первым.

Поэтому $\mathbf{y}^{(1)} \approx \frac{c_j}{\lambda_j - \lambda_j^*} e_j$ и вектор $\mathbf{y}^{(1)}$ оказывается близким по направлению к собственному вектору e_j .

Можно показать, что вектор $\mathbf{x}^{(k)}$, вычисляемый на k -й итерации, имеет вид

$$\mathbf{x}^{(k)} = \beta^{(k)} \left[c_j e_j + \sum_{\substack{i=1 \\ i \neq j}}^m \left(\frac{\lambda_j - \lambda_j^*}{\lambda_i - \lambda_j^*} \right)^k c_i e_i \right],$$

где $|\beta^{(k)}| \rightarrow |c_j^{-1}|$ при $k \rightarrow \infty$. Вектор $\mathbf{x}^{(k)}$ сходится к e_j по направлению со скоростью геометрической прогрессии, знаменатель которой

$$q = \max_{i \neq j} \frac{|\lambda_j - \lambda_j^*|}{|\lambda_i - \lambda_j^*|}.$$

Если абсолютная погрешность значения λ_j^* много меньше расстояния от λ_j до ближайшего из остальных собственных чисел (что эквивалентно выполнению условия $|\lambda_j - \lambda_j^*| \ll |\lambda_i - \lambda_j^*|$ для всех $i \neq j$), то метод обратных итераций сходится очень быстро. Чаще всего достаточно сделать 1—3 итерации.

Пример 8.6. Используя метод обратных итераций, найдем на 6-разрядном десятичном компьютере собственный вектор матрицы (8.4), отвечающий собственному значению $\lambda_1 \approx \lambda_1^* = -7.8728$.

Возьмем $\mathbf{x}^{(0)} = (1, 1, 1)^T$. Тогда система (8.26) при $k = 0$ примет вид

$$\begin{aligned} 9.8728y_1 - 9y_2 + 5y_3 &= 1, \\ 1.2y_1 + 2.4729y_2 + 6y_3 &= 1, \\ y_1 - y_2 + 0.3728y_3 &= 1. \end{aligned}$$

Вычисляя ее решение методом Гаусса, получаем $y_1^{(1)} = -123909$, $y_2^{(1)} = 99398.3$, $y_3^{(1)} = 65749$. После нормировки имеем $x_1^{(1)} = -0.720736$, $x_2^{(1)} = -0.578166$, $x_3^{(1)} = 0.382440$.

Подставляя в правую часть системы (8.26) вектор $\mathbf{x}^{(1)}$ и вычисляя решение, находим $y_1^{(2)} = -59671.5$, $y_2^{(2)} = -47867.7$, $y_3^{(2)} = 31663.1$. Полученные после нормировки значения $x_1^{(2)} = -0.720737$, $x_2^{(2)} = -0.578166$, $x_3^{(2)} = 0.382440$ с машинной точностью совпадают со значениями, полученными на первой итерации. Итерации следует прервать и считать результатом вектор $e_1 \approx (-0.720737, -0.578166, 0.382440)^T$. ▲

Замечание 1. Так как λ_j^* почти совпадает со значением λ_j , при котором $\det(\mathbf{A} - \lambda_j \mathbf{E}) = 0$, то матрица $\mathbf{A} - \lambda_j^* \mathbf{E}$ очень плохо обусловлена. Возникает естественное опасение, что большие погрешности, неизбежные при реализации вычислительного процесса (8.26), (8.27) на компьютере, могут существенно повлиять на свойства приближений. К счастью, это не так; погрешность, возникающая при численном решении системы (8.26) (которая может быть сравнима по величине с $y^{(k)}$), оказывается почти пропорциональной вектору e_j . В данном случае плохая обусловленность системы не ухудшает, а улучшает ситуацию.

В справедливости сказанного легко убедиться, если повторить вычисления из примера 8.5, используя калькулятор. Полученные значения $\mathbf{y}^{(1)}$ и $\mathbf{y}^{(2)}$ наверняка будут иметь мало общего с указанными в примере. Тем не менее приближения $\mathbf{x}^{(1)}$ и $\mathbf{x}^{(2)}$ практически совпадут с вычисленными выше, отличаясь, возможно, только знаком.

Замечание 2. Записывая равенство (8.26) в виде $\mathbf{y}^{(k+1)} = (\mathbf{A} - \lambda_j^* \mathbf{E})^{-1} \mathbf{x}^{(k)}$, замечаем, что метод обратных итераций — это просто степенной метод, примененный к матрице $(\mathbf{A} - \lambda_j^* \mathbf{E})^{-1}$.

2. Метод обратных итераций с использованием отношения Рэлея. Одной из проблем применения метода обратных итераций является необходимость получения хорошего приближения к собственному значению. Когда \mathbf{A} — симметричная матрица, можно попы-

таться использовать для оценки λ , отношение Рэлея. Этот подход приводит к следующему комбинированному методу:

$$\lambda^{(k+1)} = \rho(\mathbf{x}^{(k)}) = (\mathbf{A}\mathbf{x}^{(k)}, \mathbf{x}^{(k)}), \quad k \geq 0; \quad (8.29)$$

$$(\mathbf{A} - \lambda^{(k+1)}\mathbf{E})\mathbf{y}^{(k+1)} = \mathbf{x}^{(k)}; \quad (8.30)$$

$$\mathbf{x}^{(k+1)} = \frac{\mathbf{y}^{(k+1)}}{\|\mathbf{y}^{(k+1)}\|_2}. \quad (8.31)$$

Предполагается, что вектор $\mathbf{x}^{(0)}$ нормирован, т.е. $\|\mathbf{x}^{(0)}\| = 1$.

Если начальное приближение $\mathbf{x}^{(0)}$ хорошо аппроксимирует по направлению вектор e_j , то метод (8.29)–(8.31) сходится очень быстро. Например, если λ_j — простое собственное значение, то сходимость оказывается кубической. Одним из способов получения удовлетворительного начального приближения является выполнение нескольких итераций степенного метода.

Замечание. В случае, когда метод (8.29)–(8.31) сходится, он позволяет одновременно эффективно вычислять и собственное значение λ_j , и соответствующий собственный вектор e_j .

§ 8.4. QR-алгоритм

В настоящее время лучшим методом вычисления всех собственных значений квадратных заполненных матриц общего вида (умеренного порядка) является *QR*-алгоритм¹.

1. Основной QR-алгоритм. Опишем итерационную процедуру, являющуюся основой алгоритма. Она существенно использует возможность разложения произвольной матрицы в произведение ортогональной и верхней треугольной матриц, т.е. так называемое *QR*-разложение (см. § 5.10).

На 1-й итерации с помощью метода отражений или метода вращений вычисляют *QR*-разложение матрицы $\mathbf{A}^{(0)} = \mathbf{A}$, имеющее вид

$$\mathbf{A}^{(0)} = \mathbf{Q}_1 \mathbf{R}_1. \quad (8.32)$$

Затем строят матрицу $\mathbf{A}^{(1)} = \mathbf{R}_1 \mathbf{Q}_1$. Заметим, что из равенства (8.32) следует, что $\mathbf{R}_1 = \mathbf{Q}_1^{-1} \mathbf{A}^{(0)}$ и поэтому $\mathbf{A}^{(1)} = \mathbf{Q}_1^{-1} \mathbf{A}^{(0)} \mathbf{Q}_1$. Таким образом, матрицы $\mathbf{A}^{(1)}$ и $\mathbf{A}^{(0)}$ подобны (см. § 8.1) и поэтому имеют общий набор собственных значений $\lambda_1, \lambda_2, \dots, \lambda_m$.

¹ Этот метод независимо был предложен в 1960 г. в России математиком В. Н. Кублановской и в 1961 г. в Англии системным программистом Дж. Фрэнсисом.

На 2-й итерации находят QR -разложение матрицы $\mathbf{A}^{(1)}$, имеющее вид $\mathbf{A}^{(1)} = \mathbf{Q}_2 \mathbf{R}_2$, и вычисляют матрицу $\mathbf{A}^{(2)} = \mathbf{R}_2 \mathbf{Q}_2$, подобную матрице $\mathbf{A}^{(1)}$.

На $(k+1)$ -й итерации вычисляют разложение $\mathbf{A}^{(k)} = \mathbf{Q}_{k+1} \mathbf{R}_{k+1}$ и строят матрицу $\mathbf{A}^{(k+1)} = \mathbf{R}_{k+1} \mathbf{Q}_{k+1}$. Неограниченное продолжение этого процесса дает последовательность матриц $\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(n)}, \dots$, подобных матрице \mathbf{A} .

Обратим внимание на то, что обычной поэлементной сходимости или сходимости по норме QR -алгоритм не гарантирует. Сходимость QR -алгоритма — это *сходимость по форме* последовательности приближений $\mathbf{A}^{(n)}$ к некоторой верхней треугольной или к верхней блочно-треугольной матрице $\tilde{\mathbf{A}}$, имеющей те же собственные значения, что и матрица \mathbf{A} . Сходимость по форме характеризуется сходимостью к нулю поддиагональных элементов или элементов поддиагональных блоков, тогда как наддиагональные элементы от итерации к итерации могут существенно меняться.

Практически всегда QR -алгоритм сходится. К сожалению, в общем случае, когда собственные значения матрицы \mathbf{A} могут быть кратными или комплексными, теория его сходимости весьма сложна для изложения. Поэтому ограничимся формулировкой лишь одного из наиболее известных условий сходимости — *критерия Уилкинсона*. Согласно этому критерию предполагается выполнение следующих двух условий:

1) матрица \mathbf{A} имеет простую структуру, причем модули всех собственных значений различны:

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_m|;$$

2) приведение матрицы \mathbf{A} к диагональному виду (8.7) осуществляется с помощью матрицы подобия P , у которой все ведущие главные миноры отличны от нуля.

При выполнении этих двух условий последовательность $\mathbf{A}^{(n)}$ сходится по форме к верхней треугольной матрице $\tilde{\mathbf{A}}$ вида

$$\tilde{\mathbf{A}} = \begin{bmatrix} \lambda_1 & \times & \times & \dots & \times \\ 0 & \lambda_2 & \times & \dots & \times \\ 0 & 0 & \lambda_3 & \dots & \times \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \lambda_m \end{bmatrix}. \quad (8.33)$$

Здесь крестиками помечены элементы, в общем случае не равные нулю.

Известно, что в рассматриваемом случае элементы $a_{ij}^{(k)}$ матриц $A^{(k)}$, стоящие ниже главной диагонали, сходятся к нулю со скоростью геометрической прогрессии, причем

$$|a_{ij}^{(k)}| \leq C \left| \frac{\lambda_i}{\lambda_j} \right|^k, \quad i > j, \quad (8.34)$$

т.е. скорость сходимости $a_{ij}^{(k)}$ к нулю определяется значением отношения λ_i к λ_j (заметим, что $\left| \frac{\lambda_i}{\lambda_j} \right| < 1$ при $i > j$).

Если условие 2) не выполнено, то сходимость по-прежнему имеет место, но собственные значения в матрице \tilde{A} уже не будут расположены в порядке убывания модулей.

В общем случае предельная матрица \tilde{A} — блочно-треугольная (или получающаяся из блочно-треугольной симметричной перестановкой строк и столбцов). Наличие комплексно-сопряженных пар $\lambda_k, \bar{\lambda}_k$ собственных значений у вещественной матрицы A не является препятствием для применения QR-алгоритма. Каждой такой паре в предельной матрице будет отвечать некоторый диагональный блок — квадратная подматрица порядка 2, собственные числа которой совпадают с $\lambda_k, \bar{\lambda}_k$.

2. Ускорение QR-алгоритма. Приведенный выше вариант QR-алгоритма очень неэффективен по двум основным причинам. Первая из них состоит в том, что для реализации только одной итерации этого метода требуется выполнить порядка m^3 арифметических операций, если A — матрица общего вида. Вторая причина заключается в том, что при наличии двух близких собственных значений $\lambda_i \approx \lambda_{i-1}$ метод сходится очень медленно, так как, например элемент $a_{i,i-1}^{(k)}$ сходится к нулю со скоростью геометрической прогрессии, знаменатель которой

$$q = \left| \frac{\lambda_i}{\lambda_{i-1}} \right| \approx 1.$$

Для того чтобы уменьшить число арифметических операций, матрицу A предварительно с помощью подобных преобразований отражения или вращения приводят к так называемой *форме Хессенберга*¹:

¹ Герхард Хессенберг (1874—1925) — немецкий математик.

$$H = \begin{bmatrix} h_{11} & h_{12} & \dots & h_{1,m-1} & h_{1m} \\ h_{21} & h_{22} & \dots & h_{2,m-1} & h_{2m} \\ 0 & h_{32} & \dots & h_{3,m-1} & h_{3m} \\ 0 & 0 & \dots & h_{m,m-1} & h_{mm} \end{bmatrix}. \quad (8.35)$$

Матрица H , в которой равны нулю все элементы h_{ij} , такие, что $i > j + 1$ (т.е. все элементы, расположенные ниже диагонали, непосредственно примыкающей снизу к главной диагонали), называется *матрицей Хессенберга*. Существуют эффективные стандартные процедуры преобразования матрицы A к виду (8.35), поэтому мы не будем останавливаться подробно на этом преобразовании. Для дальнейшего важно то, что матрицы A и H подобны и обладают общим набором собственных значений, а матрица подобия P , для которой выполняется равенство

$$H = P^{-1}AP, \quad (8.36)$$

ортогональна.

После преобразования матрицы A к виду (8.35) к матрице H применяют *QR*-алгоритм. Эффективность такого подхода обусловлена наличием следующих двух замечательных свойств матриц Хессенберга.

1°. *Матрицы $H^{(k)}$, порождаемые QR-алгоритмом из матрицы $H^{(0)} = H$, сами являются матрицами Хессенберга, т.е. для них $h_{ij}^{(k)} = 0$ при $i > j + 1$.*

2°. *Для выполнения одной итерации QR-алгоритма для матрицы Хессенберга требуется число арифметических операций, пропорциональное m^2 .*

Однако, как уже было отмечено, даже достигнутая благодаря предварительному преобразованию матрицы A к виду (8.35) существенная экономия числа арифметических операций недостаточна для практического использования *QR*-алгоритма. Дело в том, что при наличии близких соседних собственных значений $\lambda_i \approx \lambda_{i-1}$ элемент $h_{i,i-1}^{(k)}$ убывает очень медленно пропорционально q^k , где $q = \left| \frac{\lambda_i}{\lambda_{i-1}} \right| \approx 1$. Для решения этой проблемы используют различные варианты *QR*-алгоритма со сдвигами.

Поясним суть этого подхода. Допустим, что для λ_i известно хорошее приближение λ_i^* . Тогда собственными значениями матрицы $\tilde{H}^{(k)} =$

$= \mathbf{H}^{(k)} - \lambda_i^* \mathbf{E}$ являются $\tilde{\lambda}_j = \lambda_j - \lambda_i^*$, $j = 1, 2, \dots, m$. В этом случае вместо отношения $\frac{\lambda_i}{\lambda_{i-1}} \approx 1$ скорость убывания поддиагонального элемента

$\tilde{h}_{i,i-1}^{(k)}$ определяет величина $\frac{\tilde{\lambda}_i}{\lambda_{i-1}} = \frac{\lambda_i - \lambda_i^*}{\lambda_{i-1} - \lambda_i^*} \approx 0$. После нескольких итераций QR-алгоритма, которые практически сделают элемент $\tilde{h}_{i,i-1}^{(k)}$ равным нулю, следует выполнить обратный сдвиг, положив $\mathbf{H}^{(k)} = \tilde{\mathbf{H}}^{(k)} + \lambda_i^* \mathbf{E}$. После выполнения этой операции матрицы \mathbf{A} и $\mathbf{H}^{(k)}$ снова имеют общий набор собственных значений.

Последовательное осуществление сдвигов для $i = m, m-1, \dots, 1$, сопровождаемых итерациями по QR-алгоритму, дает возможность быстро привести матрицу \mathbf{A} к виду (8.33). Остающийся невыясненным вопрос о том, как получить приближенные значения $\lambda_i^* \approx \lambda_i$, снимается, если учесть, что в ходе QR-алгоритма диагональные элементы $h_{ii}^{(k)}$ сходятся к λ_i при $k \rightarrow \infty$. Следовательно, в качестве λ_i^* можно, например, брать элементы¹ $h_{ii}^{(k)}$.

Итак, прежде чем применять QR-алгоритм, следует преобразовать исходную матрицу \mathbf{A} к форме Хессенберга. Без такого преобразования QR-алгоритм практически не применяется. Затем целесообразно использовать один из вариантов QR-алгоритма со сдвигами.

Пусть теперь собственные значения найдены и требуется найти один собственный вектор e_j матрицы \mathbf{A} , отвечающий собственному значению λ_j , или несколько собственных векторов. Тогда целесообразно сначала найти соответствующий собственный вектор v_j матрицы \mathbf{H} (например, методом обратных итераций), а затем вычислить e_j по формуле $e_j = P \cdot v_j$, где P — матрица подобия из (8.36).

Замечание 1. Вычисление собственного вектора e_j непосредственным применением метода обратных итераций к матрице \mathbf{A} возможно, но потребует большего числа арифметических операций по сравнению с указанным выше подходом.

¹ Чаще всего в библиотечных программах используется либо указанная стратегия сдвигов (*сдвигов по Рэлею*), либо стратегия сдвигов по Уилкинсону [103].

Замечание 2 Проведенное в этом параграфе обсуждение алгоритма носило в значительной степени ознакомительный характер. Практически не был затронут случай, когда матрица имеет кратные или комплексные собственные значения. Не рассматривались и особенности применения *QR*-алгоритма для комплексных матриц.

Замечание 3 *QR*-алгоритм обладает хорошей обусловленностью. Например, как показано в [23], в одном из его вариантов после числа итерации, не превосходящего 5, для каждого собственного значения получаются приближения $\lambda_1^*, \lambda_2^*, \dots, \lambda_m^*$, являющиеся точными собственными значениями некоторой матрицы A_* , такой, что

$$\|A - A_*\|_E \lesssim 30m^2 \varepsilon_M \|A\|_E,$$

(это утверждение сформулировано в терминах обратного анализа ошибок).

§ 8.5. Дополнительные замечания

1. В данной книге не рассмотрены некоторые весьма популярные методы решения проблемы собственных значений. Изложение *метода бисекций*, *метода вращении Якоби*, *QL-алгоритма* (являющегося вариантом *QR*-алгоритма), *LR-алгоритма* и других методов можно найти, например в [23, 24, 31, 75, 102, 103].

Авторы советуют обратить внимание на книги [49, 25*], содержащие изложение современных численных методов решения проблемы собственных значений (в том числе, подробное обсуждение *QR*-алгоритма), вполне доступное для студента или выпускника технического вуза.

2. Если A — заполненная матрица общего вида умеренного порядка, то лучшим выбором для вычисления всех собственных значений служит один из вариантов *QR*-алгоритма со сдвигами. Необходимо только предварительно преобразовать матрицу к форме Хессенберга. Часто до этого производится уравновешивание (или масштабирование) матрицы [49].

3. Когда A — симметричная матрица умеренного порядка, ее обычно приводят сначала с помощью последовательных преобразований Хаусхолдера к трехдиагональному виду. Для вычисления собственных значений полученной трехдиагональной матрицы можно использовать *QR*-алгоритм, но, по-видимому, чаще более предпочтительным является метод бисекций. Одно из достоинств этого алгоритма, состоит в том, что он позволяет находить не все собственные значения, а одно или группу нужных собственных значений.

4 Если приближенное значение собственного числа найдено, то подходящим методом вычисления соответствующего собственного вектора является метод обратных итераций

5 Методы, которые используются в настоящее время для решения проблемы собственных значений, когда A — разреженная матрица большой размерности, можно разделить на две основные группы *методы одновременных итераций* (или *итерирований подпространства*) и *методы типа Ланцоша*. Их обсуждение можно найти, например в [31, 49]. Один из простейших возможных подходов — степенной метод — был рассмотрен в § 8.2

МЕТОДЫ ОДНОМЕРНОЙ МИНИМИЗАЦИИ

Одно из важнейших направлений в конструировании изделий, а также проектировании и эксплуатации технологических процессов состоит в оптимизации (минимизации или максимизации) некоторой характеристики $f(x)$. Функцию $f(x)$ часто называют *целевой функцией*. Заметим, что основное внимание может быть уделено минимизации целевой функции, так как максимизация сводится к минимизации с помощью введения новой целевой функции $\tilde{f}(x) = -f(x)$. Когда варьируется один скалярный параметр x , возникает задача одномерной минимизации.

Необходимость изучения методов решения задачи одномерной минимизации определяется не только тем, что задача может иметь самостоятельное значение, но и в значительной мере тем, что алгоритмы одномерной минимизации являются существенной составной частью алгоритмов решения задач многомерной минимизации (см. гл. 10), а также других вычислительных задач.

§ 9.1. Задача одномерной минимизации

1. Постановка задачи. Определения. Пусть $f(x)$ — действительная функция одной переменной, определенная на множестве $X \subset (-\infty, \infty)$. Напомним, что точка $\bar{x} \in X$ называется точкой глобального минимума функции f на множестве X , если для всех $x \in X$ выполняется неравенство $f(\bar{x}) \leq f(x)$. В этом случае значение $f(\bar{x})$ называется минимальным значением функции f на X .

Точка $\bar{x} \in X$ называется точкой локального минимума функции f , если существует такая δ -окрестность этой точки, что для всех x из множества X , содержащихся в указанной δ -окрестности, выполняется неравенство $f(\bar{x}) \leq f(x)$. Если же для всех таких x , не совпадающих с \bar{x} выполняется неравенство $f(\bar{x}) < f(x)$, то \bar{x} называется точкой строгого локального минимума.

Пример 9.1. Для функции, график которой изображен на рис. 9.1, точки \bar{x}_3 и \bar{x}_4 являются точками строго локального минимума, а в точках x , удовлетворяющих неравенству $\bar{x}_1 \leq x \leq \bar{x}_2$, реализуется нестрогий локальный минимум. ▲

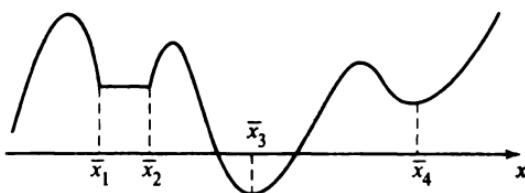


Рис. 9.1

Известно, что необходимым условием того, чтобы внутренняя для множества X точка \bar{x} была точкой локального минимума дифференцируемой функции f , является выполнение равенства

$$f'(\bar{x}) = 0. \quad (9.1)$$

Число \bar{x} , удовлетворяющее этому равенству, называется *стационарной точкой функции* f . Конечно, не всякая стационарная точка \bar{x} обязана быть точкой локального минимума. Для дважды непрерывно дифференцируемой функции достаточным условием того, чтобы стационарная точка \bar{x} была точкой строгого локального минимума, является выполнение неравенства $f''(\bar{x}) > 0$.

Существуют различные постановки задачи минимизации. В самой широкой постановке требуется найти все точки локального минимума и отвечающие им значения функции f . В приложениях чаще всего возникает задача вычисления конкретной точки локального минимума или точки глобального минимума. Иногда представляет интерес только лишь минимальное значение целевой функции, независимо от того, в какой именно точке оно достигается.

2. Отрезок локализации. Подобно тому, как алгоритмы решения нелинейных уравнений настроены на отыскание одного изолированного корня (см. гл. 4), большинство алгоритмов минимизации осуществляет лишь поиск точки локального минимума функции f . Для того чтобы применить один из таких алгоритмов минимизации, следует предварительно найти содержащий точку \bar{x} отрезок $[a, b]$, на котором она является единственной точкой локального минимума. Этот отрезок в дальнейшем будем называть *отрезком локализации*¹ точки \bar{x} . К сожалению, не существует каких-либо общих рецептов относительно того, как найти отрезок локализации. В одномерном случае полезным может оказаться табулиро-

¹ В теории оптимизации отрезок $[a, b]$ чаще называют интервалом неопределенности. Мы понимаем интервал неопределенности иначе (см. § 9.2).

вание функции с достаточно мелким шагом и (или) построение графика. Отрезок $[a, b]$ может быть известен из физических соображений, из опыта решения аналогичных задач и т.д. Для некоторых алгоритмов (например, для метода Ньютона) достаточно иметь не отрезок локализации, а хорошее начальное приближение $x^{(0)}$ к \bar{x} .

Пример 9.2. Для функции $f(x) = x^3 - x + e^{-x}$, произведем локализацию точек локального минимума.

Из графика функции, изображенного на рис. 9.2, видно, что функция $f(x)$ имеет две точки локального минимума \bar{x}_1 и \bar{x}_2 , первая из которых является и точкой глобального минимума. Для точки \bar{x}_1 за отрезок локализации можно принять отрезок $[-4, -3]$, а для точки \bar{x}_2 — отрезок $[0, 1]$.

Докажем теперь, что на отрезке $[0, 1]$ действительно содержится точка локального минимума. Для этого заметим, что $f'(x) = 3x^2 - 1 - e^{-x}$ и $f'(0) = -2 < 0$, $f'(1) = 2 - e^{-1} > 0$. Так как значения $f'(0)$ и $f'(1)$ имеют разные знаки, то на отрезке $[0, 1]$ содержится точка \bar{x} , для которой $f'(\bar{x}) = 0$. Но $f''(x) = 6x + e^{-x} > 0$ для всех $x \in [0, 1]$. Следовательно,

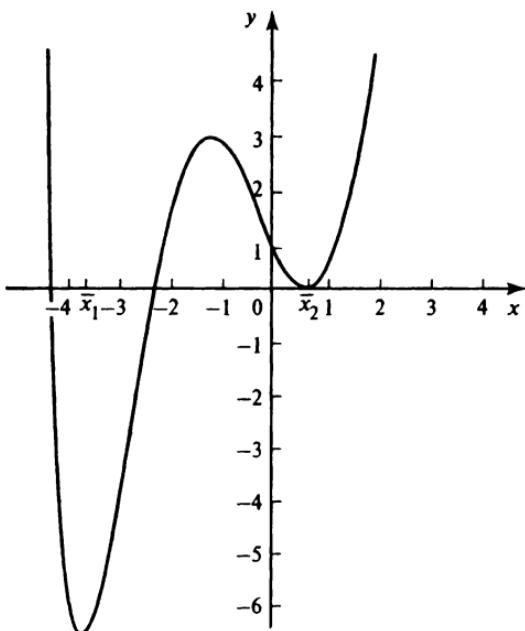


Рис. 9.2

$f''(\bar{x}) > 0$ и точка \bar{x} на отрезке $[0, 1]$ есть единственная точка локального минимума. Аналогично доказывается, что отрезок $[-4, -3]$ также является отрезком локализации. ▲

3. Унимодальные функции. Пусть f — функция, определенная на отрезке $[a, b]$. Предположим, что на этом отрезке содержится единственная точка \bar{x} локального минимума функции f , причем функция строго убывает при $x \leq \bar{x}$ и строго возрастает при $x \geq \bar{x}$. Такая функция называется *унимодальной*¹. Возможны три случая расположения точки \bar{x} на отрезке $[a, b]$: точка \bar{x} является внутренней для отрезка, \bar{x} совпадает с левым концом отрезка и \bar{x} совпадает с правым концом отрезка. Соответственно и график унимодальной функции может иметь одну из форм, схематично изображенных на рис. 9.3.

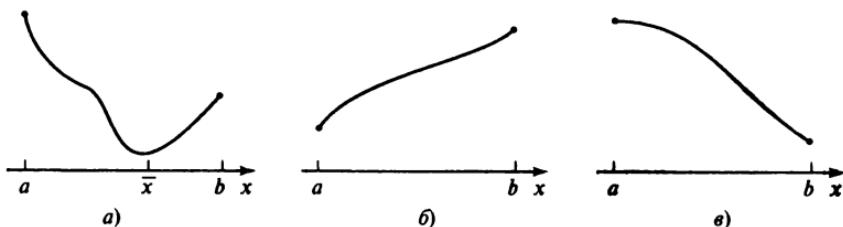


Рис. 9.3

Замечание. Унимодальная функция, вообще говоря, не обязана быть непрерывной. Например, функция, изображенная на рис. 9.4, унимодальна и имеет три точки разрыва.

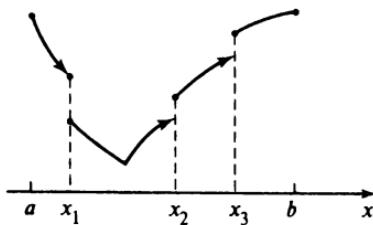


Рис. 9.4

¹ Иногда такую функцию называют строго унимодальной, а унимодальной называют функцию, которая строго убывает при $x \leq \bar{x}_1$, равна постоянной при $\bar{x}_1 \leq x \leq \bar{x}_2$ и строго возрастает при $x \geq \bar{x}_2$ [20].

Приведем достаточное условие унимодальности функции на отрезке $[a, b]$.

Предложение 9.1. Если для всех $x \in [a, b]$ выполнено условие $f''(x) > 0$, то функция унимодальна на отрезке $[a, b]$. ■

Пример 9.3. Функция $f(x) = x^3 - x + e^{-x}$ унимодальна на каждом из отрезков $[-4, -3]$ и $[0, 1]$. Чтобы убедиться в этом, достаточно заметить, что $f''(x) = 6x + e^{-x} > 0$ для всех $x \in [-4, -3]$, $x \in [0, 1]$, и воспользоваться предложением 9.1. ▲

Для сужения отрезка локализации точки минимума унимодальной функции полезно использовать следующее утверждение.

Предложение 9.2. Пусть f — унимодальная на отрезке $[a, b]$ функция и $a \leq \alpha < \gamma < \beta \leq b$. Тогда:

- 1° если $f(\alpha) \leq f(\beta)$, то $\bar{x} \in [\alpha, \beta]$;
- 2° если $f(\alpha) \geq f(\beta)$, то $\bar{x} \in [\alpha, \beta]$;
- 3° если $f(\alpha) \geq f(\gamma)$ и $f(\gamma) \leq f(\beta)$, то $\bar{x} \in [\alpha, \beta]$.

Доказательство. 1°. Предположим противное: $\bar{x} > \beta$. Тогда вследствие унимодальности f получим $f(\alpha) > f(\beta)$, что противоречит условию.

2°. Предположим противное: $\bar{x} < \alpha$. Тогда вследствие унимодальности f получим $f(\alpha) < f(\beta)$, что противоречит условию.

3°. В силу п. 1° имеем $\bar{x} \in [\alpha, \beta]$, а в силу п. 2° имеем $\bar{x} \in [\alpha, \beta]$. Следовательно, $\bar{x} \in [\alpha, \beta]$.

Геометрическая иллюстрация пп. 1° и 2° приведена на рис. 9.5. ■

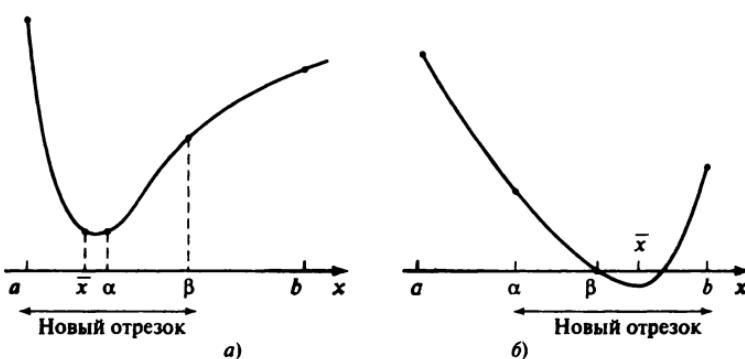


Рис. 9.5

Многие алгоритмы одномерной минимизации построены в расчете на то, что на отрезке локализации целевая функция унимодальна. В частности, такими являются алгоритмы, рассматриваемые в § 9.3.

4. Об одном подходе к локализации точки минимума. На практике часто бывает неизвестно, является ли данная функция унимодальной. Однако во многих случаях из дополнительных соображений следует, что при $x \geq x_0$ функция f сначала убывает, а затем, начиная с некоторого значения $x = \bar{x}$, становится возрастающей (правда, не исключено, что далее она снова может стать убывающей). Для того чтобы в таком случае локализовать точку \bar{x} , используют различные нестрогие методы. Один из распространенных подходов состоит в следующем. Выбирают начальный шаг $h > 0$, в несколько раз меньший предполагаемого расстояния от точки x_0 до точки \bar{x} . Затем вычисляют и сравнивают значения $f(x_0)$ и $f(x_1)$, где $x_1 = x_0 + h$.

Если оказывается, что $f(x_0) > f(x_1)$, то последовательно вычисляют значения функции f в точках $x_k = x_0 + 2^{k-1}h$ для $k \geq 2$. После обнаружения первой же точки, для которой $f(x_k) \leq f(x_{k+1})$, за отрезок локализации принимают отрезок $[x_{k-1}, x_{k+1}]$. В случае, изображенном на рис. 9.6, *a*, за отрезок локализации принят отрезок $[x_2, x_4]$.

Если же $f(x_0) \leq f(x_1)$, то последовательно вычисляют значения в точках $x_k = x_0 + h/2^{k-1}$, $k \geq 2$. После обнаружения первой же точки x_k для которой $f(x_k) < f(x_0)$, за отрезок локализации принимают отрезок

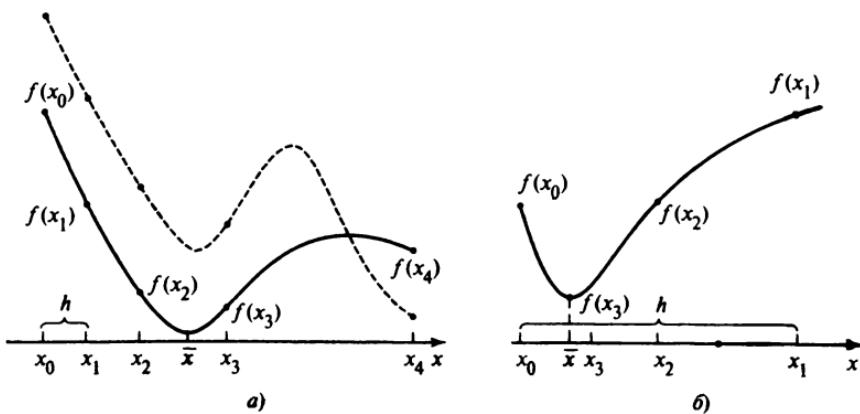


Рис. 9.6

Таблица 9.1

k	0	1	2	3	4	5
x_k	-5	-4.8	-4.6	-4.2	-3.4	-1.8
$f(x_k)$	28.4	15.7	6.74	-3.30	-5.94	2.02

$[x_0, x_{k-1}]$. В случае, изображенном на рис. 9.6, б, за отрезок локализации принят отрезок $[x_0, x_2]$.

Описанный метод не является строгим и не гарантирует, что отрезок локализации всегда будет найден. Например, для функции, график которой изображен штрихами на рис. 9.6, а, при выбранном шаге h справедливы неравенства $f(x_0) > f(x_1) > f(x_2) > f(x_3) > f(x_4)$ и поэтому отрезок локализации точки \bar{x} обнаружен уже не будет. Тем не менее этот или близкий к нему методы часто используются на практике.

Пример 9.4. Локализуем указанным выше образом точку локального минимума функции $f(x) = x^3 - x + e^{-x}$.

Возьмем $x_0 = -5$, $h = 0.2$ и положим $x_1 = x_0 + h = -4.8$. Так как $f(x_0) \approx 28.4 > f(x_1) \approx 15.7$, то будем последовательно вычислять значения функции f в точках $x_k = x_0 + 2^{k-1}h$. Из табл. 9.1 видно, что при $k = 4$ впервые выполняется неравенство $f(x_k) < f(x_{k+1})$. Поэтому за отрезок локализации следует принять отрезок $[x_3, x_5] = [-4.2, -1.8]$. ▲

§ 9.2. Обусловленность задачи минимизации

Пусть \bar{x} — точка строгого локального минимума функции f , вычисляемой с погрешностью. Будем считать, что в некоторой окрестности точки \bar{x} вычисляемые приближенные значения $f^*(x)$ удовлетворяют неравенству $|f(x) - f^*(x)| \leq \bar{\Delta} = \bar{\Delta}(f^*)$, т.е. имеют границу абсолютной погрешности, равную $\bar{\Delta}$. Как нетрудно понять, существует такая малая окрестность $(\bar{x} - \bar{\varepsilon}, \bar{x} + \bar{\varepsilon})$ точки минимума \bar{x} , для которой, основываясь на сравнении вычисляемых значений $f^*(x)$, нельзя достоверно определить ту точку, в которой действительно достигается минимум функции f . Эта ситуация схематично изображена на рис. 9.7. Интервал $(\bar{x} - \bar{\varepsilon}, \bar{x} + \bar{\varepsilon})$ будем называть *интервалом неопределенности* точки \bar{x} локального минимума.

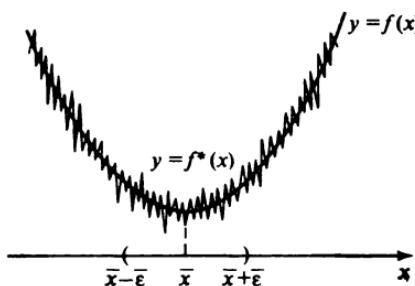


Рис. 9.7

Оценим значение $\bar{\varepsilon}$ радиуса интервала неопределенности в предположении, что функция f дважды непрерывно дифференцируема и выполнено условие $f''(\bar{x}) > 0$. В этом случае с учетом того, что $f'(\bar{x}) = 0$, для значений функции f в точках x , близких к \bar{x} , справедливо приближенное равенство

$$f(x) \approx f(\bar{x}) + \frac{f''(\bar{x})}{2} (x - \bar{x})^2.$$

Оценим минимальное расстояние между точками x и \bar{x} , начиная с которого заведомо будет выполнено неравенство $f^*(x) > f^*(\bar{x})$, т.е. точка x перестанет попадать в интервал неопределенности. Имеем

$$\begin{aligned} f^*(x) - f^*(\bar{x}) &= f(x) - f(\bar{x}) + (f^*(x) - f(x)) - (f^*(\bar{x}) - f(\bar{x})) \geq \\ &\geq f(x) - f(\bar{x}) - 2\bar{\Delta} \approx \frac{f''(\bar{x})}{2} (x - \bar{x})^2 - 2\bar{\Delta}. \end{aligned}$$

Следовательно,

$$f^*(x) - f^*(\bar{x}) \gtrsim \frac{f''(\bar{x})}{2} (x - \bar{x})^2 - 2\bar{\Delta}$$

и неравенство $f^*(x) > f^*(\bar{x})$ выполнено, если $(x - \bar{x})^2 \gtrsim 4\bar{\Delta}/f''(\bar{x})$. Таким образом,

$$\bar{\varepsilon} \approx 2\sqrt{\bar{\Delta}/f''(\bar{x})}. \quad (9.2)$$

Заметим, что любое приближение \bar{x}^* к \bar{x} , попавшее в интервал неопределенности, нельзя отличить от точного значения \bar{x} точки минимума, используя только вычисляемые значения f^* функции f . Поэтому

$$\bar{\Delta}(\bar{x}^*) \approx 2\sqrt{\bar{\Delta}(f^*)/f''(\bar{x})}. \quad (9.3)$$

Итак, рассматриваемую задачу минимизации нельзя назвать хорошо обусловленной. Если задача хорошо масштабирована, т.е. $\bar{x} \sim 1$, $|f(\bar{x})| \sim 1$, $f'(\bar{x}) \sim 1$, то соотношение (9.3) можно записать в терминах относительных погрешностей так:

$$\bar{\delta}(\bar{x}^*) \sim \sqrt{\bar{\delta}(f^*)}.$$

Отсюда следует, что если $\bar{\delta}(f^*) \sim 10^{-m}$, то $\bar{\delta}(\bar{x}^*) \sim 10^{-m/2}$. Иными словами, если значения функции вычисляются с m верными значащими цифрами, то приближенное значение точки минимума можно найти примерно лишь с $m/2$ верными значащими цифрами.

Таким образом, точность определения положения точки минимума гладкой функции существенным образом зависит от точности вычисления значений функции f . При этом если для поиска \bar{x} используются только приближенные значения $f^*(x)$, вычисляемые для различных x , то неизбежна потеря примерно половины верных значащих цифр.

Предположим теперь, что для отыскания точки локального минимума можно использовать вычисляемые каким-либо образом приближенные значения $(f')^*(x)$ производной функции f . Как уже отмечалось в § 9.1, в рассматриваемом случае задача минимизации эквивалентна задаче отыскания корня \bar{x} нелинейного уравнения $f'(x) = 0$. Из результатов § 4.2 вытекает, что последняя задача обладает значительно меньшей чувствительностью к погрешностям. В частности, справедлива следующая оценка границы абсолютной погрешности:

$$\bar{\Delta}(\bar{x}^*) \approx \frac{1}{f''(\bar{x})} \bar{\Delta}((f')^*). \quad (9.4)$$

Сравнение (9.4) с оценкой (9.3) показывает, что алгоритмы, использующие для отыскания решения \bar{x} уравнения (9.1) вычисление значений производной, могут достигать более высокой точности, чем алгоритмы, использующие для минимизации функции f только вычисление ее значений.

Пример 9.5. Оценим радиус интервала неопределенности для каждой из точек $\bar{x}_1 \approx -3.7$, $\bar{x}_2 \approx 0.7$ локального минимума функции $f(x) = x^3 - x + e^{-x}$ в случае, когда вычисление функции производится на 6-разрядном десятичном компьютере¹.

¹ Напомним, что 6-разрядным десятичным компьютером мы условились называть гипотетический компьютер, имеющий шесть десятичных разрядов мантиссы и производящий округление по дополнению.

Заметим, что $f(\bar{x}_1) \approx f(-3.7) \approx -6.5$, $f(\bar{x}_2) \approx f(0.7) \approx 0.14$. Так как для используемого компьютера¹ $\epsilon_m = 5 \cdot 10^{-7}$, то в малой окрестности точки \bar{x}_1 верхняя граница $\bar{\Delta}_1$ абсолютной погрешности вычисления f приближенно равна

$$\epsilon_m |f(\bar{x}_1)| \approx 5 \cdot 10^{-7} \cdot 6.5 = 3.25 \cdot 10^{-6}.$$

Аналогично,

$$\bar{\Delta}_2 \approx \epsilon_m |f(\bar{x}_2)| \approx 5 \cdot 10^{-7} \cdot 0.14 = 7 \cdot 10^{-8}.$$

Вычисляя значения второй производной $f''(\bar{x}) = 6x + e^{-x}$ при $x = \bar{x}_1$, $x = \bar{x}_2$, получаем $f''(\bar{x}_1) \approx 18$, $f''(\bar{x}_2) \approx 4.7$. В силу формулы (9.2) радиусы интервалов неопределенности оцениваются следующим образом:

$$\bar{\varepsilon}_1 \approx 2\sqrt{\bar{\Delta}_1/f''(\bar{x}_1)} \approx 2\sqrt{3.25 \cdot 10^{-6}/18} \approx 8 \cdot 10^{-4},$$

$$\bar{\varepsilon}_2 \approx 2\sqrt{\bar{\Delta}_2/f''(\bar{x}_2)} \approx 2\sqrt{7 \cdot 10^{-8}/4.7} \approx 2 \cdot 10^{-4}.$$

Следовательно, точку \bar{x}_2 можно найти с большей точностью, чем точку \bar{x}_1 , если использовать сравнение вычисляемых на 6-разрядном десятичном компьютере значений функции f . При этом каждую из точек можно найти с точностью $\epsilon = 10^{-3}$, но вряд ли удастся найти с точностью $\epsilon = 10^{-4}$. ▲

Пример 9.6. Пусть теперь точки \bar{x}_1 и \bar{x}_2 локального минимума функции $f(x) = x^3 - x + e^{-x}$ ищутся как решения нелинейного уравнения

$$f'(x) = 3x^2 - 1 - e^{-x} = 0. \quad (9.5)$$

Оценим радиус интервала неопределенности для каждой из точек \bar{x}_1 , \bar{x}_2 , если вычисления ведутся на том же компьютере, что и в примере 9.5.

Оценим сначала границу абсолютной погрешности вычисления производной исходя из приближенного равенства

$$\bar{\Delta} = \bar{\Delta}(f') = \bar{\Delta}(3\bar{x}^2) + \bar{\Delta}(1) + \bar{\Delta}(e^{-\bar{x}}) \approx \epsilon_m (3\bar{x}^2 + e^{-\bar{x}}).$$

Тогда

$$\bar{\Delta}_1 \approx \epsilon_m (3\bar{x}_1^2 + e^{-\bar{x}_1}) \approx 4 \cdot 10^{-5}, \quad \bar{\Delta}_2 \approx \epsilon_m (3\bar{x}_2^2 + e^{-\bar{x}_2}) \approx 10^{-6}.$$

¹ Напомним, что через ϵ_m обозначено машинное эпсилон — величина, характеризующая относительную точность представления чисел в компьютере.

На основании формулы (9.4) имеем

$$\bar{\varepsilon}_1 \approx \bar{\Delta}_1 / f''(\bar{x}_1) \approx 2 \cdot 10^{-6}, \quad \bar{\varepsilon}_2 \approx \bar{\Delta}_2 / f''(\bar{x}_2) \approx 2 \cdot 10^{-7}. \quad (9.6)$$

Заметим, что погрешности представления чисел \bar{x}_1 , \bar{x}_2 на 6-разрядном десятичном компьютере таковы: $\varepsilon_m |\bar{x}_1| \approx 2 \cdot 10^{-6}$, $\varepsilon_m |\bar{x}_2| \approx 4 \cdot 10^{-7}$. Поэтому полученные оценки (9.6) означают, что, решая уравнение (9.5), можно найти точки \bar{x}_1 и \bar{x}_2 с максимальной для используемого компьютера точностью, равной соответственно $2 \cdot 10^{-6}$ и $4 \cdot 10^{-7}$ (ср. с результатом примера 9.5). ▲

§ 9.3. Методы прямого поиска.

Оптимальный пассивный поиск.

Метод деления отрезка пополам.

Методы Фибоначчи и золотого сечения

Ряд методов минимизации основан на сравнении значений функции f , вычисляемых в точках x_1, x_2, \dots, x_N . Эти методы часто называют *методами прямого поиска*, а точки x_i — *пробными точками*.

Прежде чем перейти к изложению некоторых из наиболее известных методов прямого поиска, уточним постановку задачи. Будем считать, что требуется найти приближение \bar{x}^* к точке минимума \bar{x} унимодальной на отрезке $[a, b]$ функции f . Предположим также, что число пробных точек N заранее фиксируется и за приближение \bar{x}^* к точке минимума принимается одна из этих точек.

1. Оптимальный пассивный поиск. Метод решения поставленной задачи, в котором задается правило вычисления сразу всех пробных точек x_1, x_2, \dots, x_N и за \bar{x}^* принимается та точка x_k , для которой $f(x_k) = \min_{1 \leq i \leq N} f(x_i)$, называется *методом пассивного поиска*. Соответствующая геометрическая иллюстрация приведена на рис. 9.8.

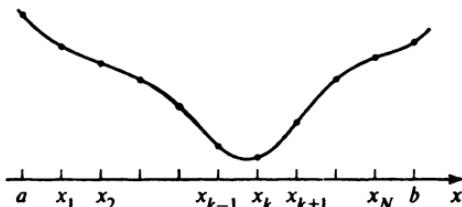


Рис. 9.8

Оценим погрешность этого метода. Для удобства положим $x_0 = a$, $x_{N+1} = b$ и будем считать, что $x_0 \leq x_1 < x_2 < \dots < x_N \leq x_{N+1}$. В силу выбора точки $\bar{x}^* = x_k$ справедливы неравенства $f(x_{k-1}) \geq f(x_k)$ и $f(x_k) \leq f(x_{k+1})$. Поэтому из п. 3° предложения 9.2 следует, что $\bar{x} \in [x_{k-1}, x_{k+1}]$. Значит

$$|\bar{x} - x_k| \leq \max\{x_k - x_{k-1}, x_{k+1} - x_k\}.$$

Так как положение точки минимума \bar{x} на отрезке $[a, b]$ заранее неизвестно, то для $\bar{x}^* = \bar{x}_k$ справедлива лишь следующая гарантированная оценка погрешности:

$$|\bar{x} - \bar{x}^*| \leq \max_{1 \leq i \leq N+1} |x_i - x_{i-1}|. \quad (9.7)$$

Можно показать, что величина, стоящая в правой части неравенства (9.7), станет минимальной, если точки x_1, x_2, \dots, x_N расположить на отрезке $[a, b]$ равномерно в соответствии с формулой $x_i = a + ih$, где $h = \Delta/(N+1)$, $\Delta = b - a$. Метод с таким выбором пробных точек называется *оптимальным пассивным поиском*. Гарантированная оценка погрешности для него выглядит так:

$$|\bar{x} - \bar{x}^*| \leq \frac{b-a}{N+1} = \frac{\Delta}{N+1}. \quad (9.8)$$

Пример 9.7. Используем оптимальный пассивный поиск для того, чтобы найти с точностью $\varepsilon = 0.1$ точку \bar{x} локального минимума функции $f(x) = x^3 - x + e^{-x}$, локализованную на отрезке $[0, 1]$.

Из (9.8) следует, что для решения задачи потребуется вычислить значения функции в девяти пробных точках вида $x_i = 0.1i$, где $i = 1, 2, \dots, 9$. Приведем таблицу этих значений (табл. 9.2).

Таблица 9.2

x	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
y	0.81	0.63	0.47	0.33	0.23	0.17	0.14	0.16	0.24

Так как минимальное значение достигается в точке $x_1 = 0.7$, то $x = 0.7 \pm 0.1$. Если бы мы попытались найти \bar{x} с точностью $\varepsilon = 10^{-2}$, то оптимальный пассивный поиск потребовал бы вычисления значений функции уже в 99 точках. ▲

2. Метод деления отрезка пополам. Пусть для решения поставленной задачи последовательно вычисляются значения функции f в N пробных точках x_1, x_2, \dots, x_N , причем для определения каждой из точек x_k можно использовать информацию о значениях функции во всех пре-

дыдущих точках x_1, x_2, \dots, x_{k-1} . Соответствующие методы называют *методами последовательного поиска*. Рассмотрим простейший из методов этого семейства — *метод деления отрезка пополам*. В нем, как и в двух других рассматриваемых в этом параграфе методах (методах Фибоначчи и золотого сечения), используется принцип последовательного сокращения отрезка локализации, основанный на предложении 9.2 и на следующем простом утверждении.

Предложение 9.3. *Если функция унимодальна на отрезке $[a, b]$, то она унимодальна и на любом отрезке $[c, d] \subset [a, b]$.* ■

Для удобства изложения обозначим отрезок $[a, b]$ через $[a^{(0)}, b^{(0)}]$. Поиск минимума начинают с выбора на отрезке $[a^{(0)}, b^{(0)}]$ двух симметрично расположенных точек

$$\alpha^{(0)} = \frac{a^{(0)} + b^{(0)}}{2} - \delta, \quad \beta^{(0)} = \frac{a^{(0)} + b^{(0)}}{2} + \delta.$$

Здесь $0 < \delta < (b - a)/2$, δ — параметр метода. Далее вычисляют значения $f(\alpha^{(0)})$ и $f(\beta^{(0)})$. Сравнение этих значений в силу предложения 9.2 позволяет сократить отрезок локализации следующим образом:

если $f(\alpha^{(0)}) \leq f(\beta^{(0)})$, то $\bar{x} \in [a^{(1)}, b^{(1)}] = [a^{(0)}, \beta^{(0)}]$;

если $f(\alpha^{(0)}) > f(\beta^{(0)})$, то $\bar{x} \in [a^{(1)}, b^{(1)}] = [\alpha^{(0)}, b^{(0)}]$.

Если описанную процедуру принять за одну итерацию метода и продолжить аналогичные операции для построения последовательности сокращающихся отрезков локализации, то получим итерационный метод. Опишем очередную его итерацию исходя из того, что отрезок локализации $[a^{(k)}, b^{(k)}]$ уже найден.

Выполняют следующие действия:

1) вычисляют

$$\alpha^{(k)} = \frac{a^{(k)} + b^{(k)}}{2} - \delta, \quad \beta^{(k)} = \frac{a^{(k)} + b^{(k)}}{2} + \delta;$$

2) находят значения $f(\alpha^{(k)})$ и $f(\beta^{(k)})$;

3) новый отрезок локализации определяют по правилу:

если $f(\alpha^{(k)}) \leq f(\beta^{(k)})$, то $[a^{(k+1)}, b^{(k+1)}] = [a^{(k)}, \beta^{(k)}]$;

если $f(\alpha^{(k)}) > f(\beta^{(k)})$, то $[a^{(k+1)}, b^{(k+1)}] = [\alpha^{(k)}, b^{(k)}]$.

В первом случае за очередное приближение к точке минимума принимают $\bar{x}^{(k+1)} = \alpha^{(k)}$, а во втором случае $\bar{x}^{(k+1)} = \beta^{(k)}$ (рис. 9.9).

Обозначим через $\Delta^{(n)} = b^{(n)} - a^{(n)}$ длину отрезка $[a^{(n)}, b^{(n)}]$. Как нетрудно заметить, справедливо равенство

$$\Delta^{(n+1)} = \frac{\Delta^{(n)}}{2} + \delta. \quad (9.9)$$

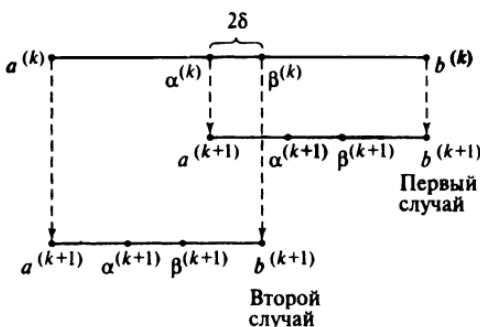


Рис. 9.9

Поэтому, если параметр δ достаточно мал ($\delta \ll \Delta^{(n)}$), то длина вновь полученного отрезка почти вдвое меньше длины предыдущего отрезка. Отсюда — и название метода.

Используя равенство (9.9), с помощью метода математической индукции легко показать, что

$$\Delta^{(n)} = \frac{(\Delta - 2\delta)}{2^n} + 2\delta.$$

Заметим, что $\Delta^{(n)}$ убывает и при $n \rightarrow \infty$ стремится к значению 2δ , оставаясь при каждом n больше него. Поэтому сделать при некотором n длину $\Delta^{(n)}$ отрезка локализации $[a^{(n)}, b^{(n)}]$ меньше заданного $\varepsilon > 0$ можно лишь, выбрав $\delta < \varepsilon/2$. В этом случае из очевидной оценки погрешности $|x^{(n)} - \bar{x}| < \Delta^{(n)}$ следует, что значение \bar{x} действительно можно найти с точностью ε и справедлив следующий критерий окончания итерационного процесса. Вычисления следует прекратить, как только окажется выполненным неравенство

$$\Delta^{(n)} \leq \varepsilon. \quad (9.10)$$

Тогда за приближение к \bar{x} с точностью ε можно принять $\bar{x}^* = x^{(n)}$.

Замечание. При реализации метода на компьютере необходимо учитывать, что вычисления значений функции f будут производиться с погрешностью. Для того чтобы знак разности $f^*(\alpha^{(n)}) - f^*(\beta^{(n)})$ совпадал со знаком разности $f(\alpha^{(n)}) - f(\beta^{(n)})$, необходимо, чтобы выполнялось условие $\delta \gtrsim \bar{\varepsilon}$ (см. равенство (9.2)). Поэтому δ нельзя задавать слишком малым.

Пример 9.8. Применив метод деления отрезка пополам, найдем с точностью $\varepsilon = 10^{-2}$ точку \bar{x} локального минимума функции $f(x) = x^3 -$

Таблица 9.3

n	$a^{(n)}$	$b^{(n)}$	$\alpha^{(n)}$	$f(\alpha^{(n)})$	$\beta^{(n)}$	$f(\beta^{(n)})$	$\Delta^{(n)}$
0	0.000000	1.000000	0.499000	0.232389	0.501000	0.230676	1.000
1	0.499000	1.000000	0.748500	0.143924	0.750500	0.144350	0.501
2	0.499000	0.750500	0.623750	0.154860	0.625750	0.154131	0.252
3	0.623750	0.750500	0.686125	0.140403	0.688125	0.140230	0.125
4	0.686125	0.750500	0.717088	0.139821	0.719088	0.139940	0.063
5	0.686125	0.719088	0.701607	0.139549	0.703607	0.139520	0.033
6	0.701607	0.719088	0.708348	0.139543	0.711348	0.139587	0.017
7	0.710141	0.719813	—	—	—	—	0.010

$-x + e^{-x}$ локализованную на отрезке $[0, 1]$. Вычисления будем вести на 6-разрядном десятичном компьютере.

Зададим $\delta = 10^{-3}$, $a^{(0)} = 0$, $b^{(0)} = 1$.

1-я итерация.

1. Вычислим:

$$\alpha^{(0)} = \frac{a^{(0)} + b^{(0)}}{2} - \delta = 0.499, \quad \beta^{(0)} = \frac{a^{(0)} + b^{(0)}}{2} + \delta = 0.501.$$

2. Определим значения $f(\alpha^{(0)}) \approx 0.232389$, $f(\beta^{(0)}) \approx 0.230676$.

3. Так как $f(\alpha^{(0)}) > f(\beta^{(0)})$, то следует положить $[a^{(1)}, b^{(1)}] = [0.499, 1]$.

Результаты следующих итераций приведены в табл. 9.3.

Так как $\Delta^{(7)} \leq \epsilon$, то при $n = 7$ итерации прекратим и положим $\bar{x} \approx \beta^{(n)} = 0.711348$. Таким образом, $\bar{x} = 0.71 \pm 0.01$. Заметим, что для достижения точности $\epsilon = 10^{-2}$ потребовалось 14 вычислений функции. ▲

3. **Метод Фибоначчи.** Заметим, что метод деления отрезка пополам требует на каждой итерации вычисления двух новых значений функции, так как найденные на предыдущей итерации в точках $\alpha^{(n)}$ и $\beta^{(n)}$ значения далее не используются. Обратим теперь внимание на то, что одна из этих точек (обозначенная в предыдущем пункте через $x^{(n)}$) является внутренней для отрезка $[a^{(n)}, b^{(n)}]$ и поэтому дальнейшее сокращение отрезка можно произвести, вычислив дополнительное значение функции лишь в одной новой точке. Это наблюдение приводит к методам, требующим на каждой итерации (кроме первой) расчета лишь одного нового значения функции f . Два наиболее известных среди них — методы Фибоначчи и золотого сечения.

Метод Фибоначчи¹ является оптимальным последовательным методом, т.е. методом, обеспечивающим максимальное гарантированное сокращение отрезка локализации при заданном числе N вычислений функции. Этот метод основан на использовании чисел Фибоначчи F_n , задаваемых рекуррентной формулой

$$F_n = F_{n-1} + F_{n-2} \quad (n \geq 2)$$

и начальными значениями $F_0 = 1$, $F_1 = 1$. Укажем несколько первых чисел: $F_0 = 1$, $F_1 = 1$, $F_2 = 2$, $F_3 = 3$, $F_4 = 5$, $F_5 = 8$, $F_6 = 13$, $F_7 = 21$, $F_8 = 34$, $F_9 = 55$, $F_{10} = 89$, $F_{11} = 144$.

Метод Фибоначчи состоит из $N - 1$ шагов. Очередной $(k + 1)$ -й шаг выполняют здесь аналогично $(k + 1)$ -й итерации метода деления отрезка пополам. В отличие от него точки $\alpha^{(k)}$, $\beta^{(k)}$ здесь находят по формулам

$$\alpha^{(k)} = a^{(k)} + \frac{F_{N-k-1}}{F_{N-k+1}} \Delta^{(k)}, \quad \beta^{(k)} = a^{(k)} + \frac{F_{N-k}}{F_{N-k+1}} \Delta^{(k)}.$$

Новый отрезок локализации определяют по тому же правилу:

если $f(\alpha^{(k)}) \leq f(\beta^{(k)})$, то $[a^{(k+1)}, b^{(k+1)}] = [\alpha^{(k)}, \beta^{(k)}]$;

если $f(\alpha^{(k)}) > f(\beta^{(k)})$, то $[a^{(k+1)}, b^{(k+1)}] = [\alpha^{(k)}, b^{(k)}]$.

В первом случае за очередное приближение к точке минимума принимают $x^{(k+1)} = \alpha^{(k)}$, а во втором случае $x^{(k+1)} = \beta^{(k)}$ (рис. 9.10).

Важно то, что в любом случае точка $x^{(k+1)}$ совпадает с одной из точек

$$\alpha^{(k+1)} = a^{(k+1)} + \frac{F_{N-k-2}}{F_{N-k}} \Delta^{(k+1)},$$

$$\beta^{(k+1)} = a^{(k+1)} + \frac{F_{N-k-1}}{F_{N-k}} \Delta^{(k+1)}.$$

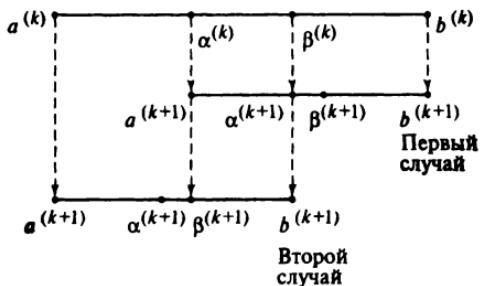


Рис. 9.10

¹ Фибоначчи (Леонардо Пизанский) (1180—1240) — итальянский математик.

Поэтому на очередном шаге достаточно вычислить значение функции лишь в одной недостающей точке.

В результате выполнения $N - 1$ шагов отрезок локализации уменьшается в $F_{N+1}/2$ раз, а точка $x^{(N-1)}$ оказывается центральной для последнего отрезка локализации $[a^{(N-1)}, b^{(N-1)}]$. Поэтому для $x^{(N-1)}$ справедлива следующая оценка погрешности:

$$|\bar{x} - x^{(N-1)}| \leq \frac{1}{F_{N+1}} \Delta. \quad (9.11)$$

Пример 9.9. Применив метод Фибоначчи, найдем с точностью $\epsilon = 10^{-2}$ точку \bar{x} локального минимума функции $f(x) = x^3 - x + e^{-x}$, локализованную на отрезке $[0, 1]$. Вычисления будем вести на 6-разрядном десятичном компьютере.

Первым среди чисел Фибоначчи, для которого выполняется условие $\Delta/F_{N+1} < \epsilon$ (где $\Delta = 1$), является число $F_{11} = 144$, отвечающее $N = 10$. Зададим $a^{(0)} = 0$, $b^{(0)} = 1$.

Первый шаг.

1. Вычислим

$$\alpha^{(0)} = a^{(0)} + (b^{(0)} - a^{(0)}) \frac{F_9}{F_{11}} = \frac{55}{144} \approx 0.381944,$$

$$\beta^{(0)} = a^{(0)} + (b^{(0)} - a^{(0)}) \frac{F_{10}}{F_{11}} = \frac{89}{144} \approx 0.618056.$$

2. Определим значения $f(\alpha^{(0)}) \approx 0.356308$, $f(\beta^{(0)}) \approx 0.157028$.

3. Так как $f(\alpha^{(0)}) > f(\beta^{(0)})$, то положим $[a^{(1)}, b^{(1)}] = [\alpha^{(0)}, b^{(0)}]$.

Второй шаг.

1. С учетом того, что $\alpha^{(1)} = \beta^{(0)} \approx 0.618056$, найдем:

$$\beta^{(1)} = a^{(1)} + (b^{(1)} - a^{(1)}) F_9 / F_{10} \approx 0.763889.$$

2. Вычислим значение $f(\beta^{(1)}) \approx 0.147712$.

3. Так как $f(\alpha^{(1)}) > f(\beta^{(1)})$, то положим $[a^{(2)}, b^{(2)}] = [\alpha^{(1)}, b^{(1)}]$.

Результаты остальных шагов приведены в табл. 9.4.

После девяти шагов вычисления прекращаем и полагаем $\bar{x} \approx \beta^{(8)} = 0.708334$. Таким образом, $\bar{x} = 0.71 \pm 0.01$. Заметим, что для достижения точности $\epsilon = 10^{-2}$ потребовалось десять вычислений значения функции, в то время как при использовании метода деления отрезка пополам необходимо 14 вычислений. ▲

Таблица 9.4

n	$a^{(n)}$	$b^{(n)}$	$\alpha^{(n)}$	$f(\alpha^{(n)})$	$\beta^{(n)}$	$f(\beta^{(n)})$	$\Delta^{(n)}$
0	0.000000	1.000000	0.381944	0.356308	0.618056	0.157028	1.000
1	0.381944	1.000000	0.618056	0.157028	0.763889	0.147712	0.618
2	0.618056	1.000000	0.763889	0.147712	0.854167	0.194672	0.382
3	0.618056	0.854167	0.708334	0.139527	0.763889	0.147712	0.236
4	0.618056	0.763889	0.673611	0.141905	0.708334	0.139527	0.119
5	0.673611	0.763889	0.708334	0.139527	0.729167	0.140830	0.090
6	0.673611	0.729167	0.694445	0.139805	0.708334	0.139527	0.056
7	0.694445	0.729167	0.708334	0.139527	0.715278	0.139731	0.035
8	0.694445	0.715278	0.701389	0.139553	0.708334	0.139527	0.021
9	0.701389	0.715278	0.708334	—	0.708334	—	0.014

Хотя метод Фибоначчи и оптимален в указанном выше смысле, часто он неудобен для использования. В частности, это касается возможного применения поиска Фибоначчи для решения одномерных подзадач в алгоритмах многомерной минимизации. Здесь нередко эффективность алгоритма одномерной минимизации оценивается не по тому, какая точность в значении \bar{x} получена, а совсем по другим критериям, к которым метод плохо приспособлен. Например, бывает важно достигнуть минимального значения функции f с некоторой относительной точностью δ либо уменьшить значение f на определенную величину.

Кроме того, число выполняемых итераций заранее фиксируется, что удобно далеко не всегда.

4. Метод золотого сечения. Из-за указанных недостатков вместо метода Фибоначчи чаще используется теоретически почти столь же эффективный *метод золотого сечения*.

Напомним, что *золотым сечением*¹ отрезка называется такое разбиение отрезка на две неравные части, что отношение длины всего отрезка к длине его большей части равно отношению длины большей части к длине меньшей части отрезка.

Золотое сечение отрезка $[a, b]$ осуществляется каждой из двух симметрично расположенных относительно центра отрезка точек

$$\alpha = a + \frac{2}{3 + \sqrt{5}} (b - a), \quad \beta = a + \frac{2}{1 + \sqrt{5}} (b - a)$$

¹ Термин «золотое сечение» ввел Леонардо да Винчи. Принципы золотого сечения широко использовались при композиционном построении многих произведений мирового искусства (в особенности, архитектурных сооружений античности и эпохи Возрождения).



(9.12)

Рис. 9.11

(рис. 9.11). Действительно, как нетрудно проверить,

$$\frac{b-a}{b-\alpha} = \frac{b-\alpha}{\alpha-a} = \frac{1+\sqrt{5}}{2}, \quad \frac{b-a}{\beta-a} = \frac{\beta-a}{b-\beta} = \frac{1+\sqrt{5}}{2}.$$

Замечательно то, что точка α осуществляет золотое сечение не только отрезка $[a, b]$, но и отрезка $[a, \beta]$. Действительно,

$$\frac{\beta-a}{\alpha-a} = \frac{\alpha-a}{\beta-\alpha} = \frac{1+\sqrt{5}}{2}.$$

Точно так же точка β осуществляет золотое сечение не только отрезка $[a, b]$, но и отрезка $[\alpha, b]$. Этот факт далее существенно используется.

Очередная $(k+1)$ -я итерация метода золотого сечения производится аналогично $(k+1)$ -й итерации метода деления отрезка пополам. В отличие от него точки $\alpha^{(k)}, \beta^{(k)}$ находятся по формулам

$$\alpha^{(k)} = a^{(k)} + \frac{2}{3+\sqrt{5}} \Delta^{(k)}, \quad \beta^{(k)} = a^{(k)} + \frac{2}{1+\sqrt{5}} \Delta^{(k)}.$$

Важно то, что какой бы из отрезков $[a^{(k)}, \beta^{(k)}]$ или $[\alpha^{(k)}, b^{(k)}]$ не был выбран за очередной отрезок локализации, точка $x^{(k+1)}$ (в первом случае $x^{(k+1)} = \alpha^{(k)}$, а во втором случае $x^{(k+1)} = \beta^{(k)}$) совпадает с одной из точек $\alpha^{(k+1)}, \beta^{(k+1)}$. Поэтому на очередном шаге достаточно вычислить значение функции лишь в одной недостающей точке.

Заметим, что точка $x^{(n)}$ отстоит от концов отрезка $[a^{(n)}, b^{(n)}]$ на величину, не превышающую $\frac{2}{1+\sqrt{5}} \Delta^{(n)}$, поэтому верна оценка

$$|x^{(n)} - \bar{x}| \leq \frac{2}{1+\sqrt{5}} \Delta^{(n)} = \Delta^{(n+1)}, \quad (9.12)$$

которую можно использовать для апостериорной оценки погрешности. Заметим, что каждая итерация сокращает длину отрезка локализации в $(1+\sqrt{5})/2$ раз. Поэтому $b^{(n)} - a^{(n)} = \Delta^{(n)} = (2/(1+\sqrt{5}))^n \Delta$ и справедлива следующая априорная оценка погрешности:

$$|x^{(n)} - \bar{x}| \leq \left[\frac{2}{1+\sqrt{5}} \right]^{n+1} \Delta. \quad (9.13)$$

Таким образом, метод золотого сечения сходится со скоростью геометрической прогрессии, знаменатель которой $q = 2/(1 + \sqrt{5}) \approx 0.62$.

Пример 9.10. Найдем методом золотого сечения с точностью $\epsilon = 10^{-2}$ точку \bar{x} локального минимума функции $f(x) = x^3 - x + e^{-x}$ локализованную на отрезке $[0, 1]$.

Положим $a^{(0)} = 0, b^{(0)} = 1$.

1-я итерация.

1. Вычислим:

$$\alpha^{(0)} = a^{(0)} + \frac{2}{3 + \sqrt{5}} (b^{(0)} - a^{(0)}) \approx 0.381966,$$

$$\beta^{(0)} = a^{(0)} + \frac{2}{1 + \sqrt{5}} (b^{(0)} - a^{(0)}) \approx 0.618034.$$

2. Найдем $f(\alpha^{(0)}) \approx 0.356280, f(\beta^{(0)}) \approx 0.157037$.

3. Так как $f(\alpha^{(0)}) > f(\beta^{(0)})$, то положим $[a^{(1)}, b^{(1)}] = [a^{(0)}, b^{(0)}]$.

2-я итерация.

1. Учтем, что $\alpha^{(1)} = \beta^{(0)} \approx 0.618034$, и вычислим:

$$\beta^{(1)} = a^{(1)} + \frac{2}{1 + \sqrt{5}} (b^{(1)} - a^{(1)}) \approx 0.763936.$$

2. Определим $f(\beta^{(1)}) \approx 0.147725$.

3. Так как $f(\alpha^{(1)}) = f(\beta^{(0)}) > f(\beta^{(1)})$, то положим $[a^{(2)}, b^{(2)}] = [\alpha^{(1)}, b^{(1)}]$.

Результаты остальных итераций приведены в табл. 9.5.

Таблица 9.5

n	$a^{(n)}$	$b^{(n)}$	$\alpha^{(n)}$	$f(\alpha^{(n)})$	$\beta^{(n)}$	$f(\beta^{(n)})$	$\Delta^{(n+1)}$
0	0.000000	1.000000	0.381966	0.356280	0.618034	0.157037	0.618
1	0.381966	1.000000	0.618034	0.157037	0.763936	0.147725	0.382
2	0.618034	1.000000	0.763936	0.147725	0.854102	0.194622	0.236
3	0.618034	0.854102	0.708204	0.139526	0.763936	0.147725	0.146
4	0.618034	0.763936	0.673764	0.141883	0.708204	0.139526	0.090
5	0.673764	0.763936	0.708204	0.139526	0.729493	0.140868	0.056
6	0.673764	0.729493	0.695051	0.139774	0.708204	0.139526	0.034
7	0.695051	0.729493	0.708204	0.139526	0.716337	0.139782	0.021
8	0.695051	0.716337	0.703182	0.139525	0.708204	0.139526	0.013
9	0.695051	0.708204	—	—	—	—	0.008

Так как $\Delta^{(10)} < \epsilon$, то итерации следует прекратить и положить $\bar{x} \approx \alpha^{(8)} \approx 0.703182$. Таким образом, $\bar{x} = 0.70 \pm 0.01$. Отметим, что для достижения точности $\epsilon = 10^{-2}$ потребовалось десять вычислений значений функции, как и в методе Фибоначчи. ▲

5. Эффективность методов прямого поиска. Эффективность указанных методов можно оценивать, например тем, во сколько раз уменьшается после использования N вычислений значений функции первоначальная длина Δ отрезка локализации. Другой критерий — значение гарантированной оценки погрешности. Данные табл. 9.6 позволяют сравнить по этим критериям рассмотренные выше методы. Как видно, очень неэффективен пассивный поиск. Метод деления отрезка пополам уступает почти эквивалентным по эффективности методам Фибоначчи и золотого сечения.

Таблица 9.6

Метод прямого поиска	Длина отрезка локализации	Гарантированная оценка погрешности
Оптимальный пассивный поиск	$\frac{2}{N+1} \Delta$	$\frac{1}{N+1} \Delta$
Метод деления отрезка пополам (N — четное, величиной δ пренебрегаем)	$\frac{1}{2^{N/2}} \Delta \approx (0.71)^N \cdot \Delta$	$\approx 0.5 \cdot (0.71)^N \cdot \Delta$
Метод Фибоначчи	$\frac{2}{F_{N+1}} \Delta \approx 1.7 \cdot (0.62)^N \cdot \Delta$	$\approx 0.85 \cdot (0.62)^N \cdot \Delta$
Метод золотого сечения	$\left[\frac{2}{1 + \sqrt{5}} \right]^{N-1} \Delta \approx 1.6 \cdot (0.62)^N \cdot \Delta$	$\approx (0.62)^N \cdot \Delta$

Приведем еще одну таблицу (табл. 9.7), из которой видно, сколько вычислений значений функции f нужно сделать для того, чтобы достигнуть точности ϵ . Предполагается, что начальный отрезок локализации имеет единичную длину.

Таблица 9.7

Метод прямого поиска	Число N при заданном ϵ , равном				
	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}
Оптимальный пассивный поиск	9	99	999	9999	99999
Метод деления отрезка пополам ($\delta \ll \epsilon$)	6	12	18	26	32
Метод Фибоначчи	5	10	15	19	24
Метод золотого сечения	5	10	15	20	24

6. Влияние погрешности вычислений. Одна из самых распространенных ошибок при обращении к стандартным программам, реали-

зующим тот или иной метод на компьютере, состоит в завышении требуемой точности. Необходимо понимать, что при наличии погрешностей в вычислении значений функции f достигимая точность ε методов прямого поиска ограничена снизу величиной $\bar{\varepsilon} \approx 2\sqrt{\Delta(f^*)/f''(\bar{x})}$, где $\bar{\varepsilon}$ — радиус интервала неопределенности (см. § 9.2). Это означает, например, что прямые методы не позволяют найти на 6-разрядном десятичном компьютере точку \bar{x}_2 локального экстремума функции $f(x) = x^3 - x + e^{-x}$ с точностью $\varepsilon < \bar{\varepsilon}_2 \approx 2 \cdot 10^{-4}$. Задание $\varepsilon < \bar{\varepsilon}_2$ приведет лишь к бесполезной трате машинного времени.

§ 9.4. Метод Ньютона и другие методы минимизации гладких функций

Отметим, что применение рассмотренных выше методов деления отрезка пополам, Фибоначчи и золотого сечения не позволяет извлечь никакой выгоды из возможной гладкости функции. Существуют методы, которые могут оказаться более эффективными, если минимизируемая функция достаточно гладкая. Часть из них является просто модификациями известных методов решения нелинейных уравнений (см. гл. 4) применительно к необходимому условию экстремума

$$f'(x) = 0. \quad (9.14)$$

1. Метод бисекции. Пусть f — унимодальная непрерывно дифференцируемая на отрезке $[a^{(0)}, b^{(0)}] = [a, b]$ функция и на отрезке $[a, b]$ точка \bar{x} является единственной стационарной точкой. Применительно к решению уравнения (9.14) одна итерация *метода бисекции* выглядит следующим образом.

Пусть отрезок локализации $[a^{(n)}, b^{(n)}]$ известен и найдено значение $x^{(n)} = (a^{(n)} + b^{(n)})/2$. Тогда производят следующие действия:

1) вычисляют значение $f'(x^{(n)})$;

2) если $f'(x^{(n)}) < 0$, то полагают $[a^{(n+1)}, b^{(n+1)}] = [x^{(n)}, b^{(n)}]$. В противном случае полагают $[a^{(n+1)}, b^{(n+1)}] = [a^{(n)}, x^{(n)}]$;

3) вычисляют $x^{(n+1)} = (a^{(n+1)} + b^{(n+1)})/2$.

В рассматриваемом случае метод сходится с оценкой погрешности

$$|x^{(n)} - \bar{x}| \leq \frac{b - a}{2^{n+1}} \quad (9.15)$$

и обладает всеми присущими методу бисекции решения нелинейных уравнений достоинствами и недостатками (см. § 4.3). Возникает лишь

дополнительная проблема, связанная с необходимостью вычисления производной f'

2. Метод Ньютона. Для решения уравнения (9.14) можно попытаться воспользоваться методом Ньютона (см. § 4.6), расчетная формула которого в данном случае принимает вид

$$x^{(n+1)} = x^{(n)} - \frac{f'(x^{(n)})}{f''(x^{(n)})}, \quad n \geq 0. \quad (9.16)$$

Следствием теоремы 4.6 является следующее утверждение

Теорема 9.1. Пусть в некоторой окрестности точки \bar{x} функция f трижды непрерывно дифференцируема и выполняется условие $f''(\bar{x}) > 0$. Тогда найдется такая малая σ -окрестность корня \bar{x} , что при произвольном выборе начального приближения $x^{(0)}$ из этой σ -окрестности метод Ньютона (9.16) сходится квадратично. ■

В силу сверхлинейной сходимости для метода Ньютона можно использовать следующий критерий окончания итераций:

$$|x^{(n)} - x^{(n-1)}| < \varepsilon. \quad (9.17)$$

Пример 9.11. Используя метод бисекции, найдем с точностью $\varepsilon = 10^{-2}$ точку локального минимума функции $f(x) = x^3 - x + e^{-x}$, локализованную на отрезке $[0, 1]$.

Положим $a^{(0)} = 0$, $b^{(0)} = 1$, $x^{(0)} = (a^{(0)} + b^{(0)})/2 = 0.5$. Заметим, что $f'(x) = 3x^2 - 1 - e^{-x}$.

1-я итерация.

- 1) вычислим $f'(x^{(0)}) \approx -0.856531$;
- 2) так как $f'(x^{(0)}) < 0$, то $[a^{(1)}, b^{(1)}] = [x^{(0)}, b^{(0)}]$;
- 3) вычислим $x^{(1)} = (a^{(1)} + b^{(1)})/2 = 0.75$.

Результаты остальных итераций приведены в табл. 9.8

Таблица 9.8

n	$a^{(n)}$	$b^{(n)}$	$x^{(n)}$	$f'(x^{(n)})$	$(b^{(n)} - a^{(n)})/2$
0	0.000000	1 000000	0.500000	-0.856531	0.500
1	0.500000	1 000000	0.750000	0.215133	0.250
2	0.500000	0 750000	0.625000	-0.363386	0.125
3	0.625000	0 750000	0.687500	-0.084863	0.063
4	0.687500	0 750000	0.718750	0.062444	0.031
5	0.687500	0 718750	0.703125	-0.011882	0.016
6	0.703125	0 718750	0.710938	—	0.008

После выполнения шести итераций вычисления можно прекратить и положить $\bar{x} \approx x^{(6)}$. Таким образом, $\bar{x} = 0.71 \pm 0.01$. Отметим, что для достижения точности $\varepsilon = 10^{-2}$ потребовалось шесть вычислений значения производной $f'(x)$. \blacktriangle

Пример 9.12. Используя метод Ньютона, найдем с точностью $\varepsilon = 10^{-6}$ точку локального минимума функции $f(x) = x^3 - x + e^{-x}$, локализованную на отрезке $[0, 1]$.

Положим $x^{(0)} = 0.5$. Результаты вычислений по формуле (9.16), имеющей в данном случае вид

$$x^{(n+1)} = x^{(n)} - \frac{3(x^{(n)})^2 - 1 - e^{-x^{(n)}}}{6x^{(n)} + e^{-x^{(n)}}},$$

приведены в табл. 9.9.

Таблица 9.9

n	$x^{(n)}$	$ x^{(n)} - x^{(n-1)} $
0	0.5000000	—
1	0.7374944	$2 \cdot 10^{-1}$
2	0.7062126	$3 \cdot 10^{-2}$
3	0.7056421	$6 \cdot 10^{-4}$
4	0.7056419	$2 \cdot 10^{-7}$

При $n = 4$ итерации прерываются. Можно считать, что $\bar{x} = 0.705642 \pm \pm 10^{-6}$. Заметим, что точность $\varepsilon = 10^{-2}$ была достигнута уже после выполнения двух итераций. \blacktriangle

3. Метод последовательной параболической интерполяции. Этот метод предназначен для минимизации гладких функций, но в отличие от методов бисекции и Ньютона не требует вычисления производных.

Опишем одну итерацию простейшего варианта этого метода. Предположим, что уже известны три предыдущих приближения $x^{(k-2)}$, $x^{(k-1)}$, $x^{(k)}$ к точке \bar{x} . Пусть $y = P_2(x) = m^{(k)} + n^{(k)}(x - x^{(k)}) + p^{(k)}(x - x^{(k)})^2$ —

уравнение параболы, проходящей через три точки плоскости с координатами $(x^{(k-2)}, f(x^{(k-2)}))$, $(x^{(k-1)}, f(x^{(k-1)}))$, $(x^{(k)}, f(x^{(k)}))$. Здесь

$$m^{(k)} = f(x^{(k)}), \quad n^{(k)} = \frac{x^{(k-1)} - x^{(k)}}{x^{(k-1)} - x^{(k-2)}} \frac{f(x^{(k)}) - f(x^{(k-2)})}{x^{(k)} - x^{(k-2)}} + \\ + \frac{x^{(k)} - x^{(k-2)}}{x^{(k-1)} - x^{(k-2)}} \frac{f(x^{(k-1)}) - f(x^{(k)})}{x^{(k-1)} - x^{(k)}}, \\ p^{(k)} = \frac{1}{x^{(k-1)} - x^{(k-2)}} \left[\frac{f(x^{(k-1)}) - f(x^{(k)})}{x^{(k-1)} - x^{(k)}} - \frac{f(x^{(k)}) - f(x^{(k-2)})}{x^{(k)} - x^{(k-2)}} \right].$$

За очередное приближение $x^{(k+1)}$ к \bar{x} принимается та точка, в которой функция $P_2(x)$ достигает минимума. Из уравнения

$$P'_2(x^{(k+1)}) = 2p^{(k)}(x^{(k+1)} - x^{(k)}) + n^{(k)} = 0$$

получается формула

$$x^{(k+1)} = x^{(k)} - \frac{n^{(k)}}{2p^{(k)}}. \quad (9.18)$$

Естественно, что для начала работы этого метода требуется выбор трех начальных приближений $x^{(0)}, x^{(1)}, x^{(2)}$.

Пусть функция f трижды непрерывно дифференцируема в некоторой окрестности точки \bar{x} и удовлетворяет условию $f''(\bar{x}) > 0$. Можно показать, что в этом случае выбор начальных приближений $x^{(0)}, x^{(1)}, x^{(2)}$ из достаточно малой окрестности точки \bar{x} гарантирует, что $p^{(k)} \neq 0$ для всех k и метод последовательной параболической интерполяции сходится сверхлинейно, с порядком, приближенно равным 1.324. Поэтому в качестве критерия окончания итерационного процесса можно принять неравенство (9.17).

Отметим, что в описанном методе используются только значения функции f , вычисляемые в точках $x^{(k)}$. Поэтому (как и для методов прямого поиска) при его реализации на компьютере достижимая точность метода ограничена снизу значением, равным радиусу $\bar{\varepsilon}$ интервала неопределенности. После того как очередное приближение $x^{(n)}$ попадет в интервал $(\bar{x} - \bar{\varepsilon}, \bar{x} + \bar{\varepsilon})$, дальнейшие вычисления теряют смысл (см. § 9.2).

Существуют различные модификации метода последовательной параболической интерполяции. Одна из них, например обеспечивает принадлежность очередного приближения предыдущему отрезку локализации и дает последовательность стягивающихся к точке \bar{x} отрезков.

4. Гибридные алгоритмы. Лучшими среди универсальных методов одномерной минимизации считаются так называемые *гибридные* (или *регуляризованные*) алгоритмы. Они представляют собой комбинации надежных, но медленно сходящихся алгоритмов типа бисекции (если возможно вычисление $f'(x)$) или золотого сечения с быстро сходящимися методами типа последовательной параболической интерполяции или Ньютона. Эти алгоритмы обладают высокой надежностью и гарантированной сходимостью, причем сходимость становится сверхлинейной, если в окрестности точки строгого минимума функция f достаточно гладкая.

Примером эффективного гибридного алгоритма является алгоритм FMIN, изложенный в [106]. Алгоритм FMIN осуществляет поиск минимума методом золотого сечения, переключаясь по возможности на параболическую интерполяцию.

§ 9.5. Дополнительные замечания

1. Дополнительную информацию о методах одномерной минимизации можно найти, например в [20].

2. Описанные выше методы приспособлены, как правило, для минимизации унимодальных функций. Если эти методы применить для минимизации непрерывной функции, не являющейся унимодальной на рассматриваемом отрезке, то мы получим, вообще говоря, лишь точку локального экстремума. Поэтому такие методы часто называют локальными методами минимизации. К настоящему времени разработан ряд методов, которые предназначены для поиска глобального минимума. С некоторыми из них можно ознакомиться в [20].

3. Решение задачи минимизации существенно усложняется, если на значения функции накладываются случайные ошибки (помехи). Так бывает, например, тогда, когда значения функции получают в результате измерений какой-либо физической величины. В том случае, когда ошибки являются случайными величинами и обладают определенными вероятностными характеристиками, для поиска минимума можно использовать *метод стохастической аппроксимации*. Понятие об этом методе можно получить из [20]; там же содержатся ссылки на соответствующую литературу.

МЕТОДЫ МНОГОМЕРНОЙ МИНИМИЗАЦИИ

Одной из наиболее часто встречающихся в прикладных расчетах и научных исследованиях вычислительных задач является задача минимизации¹ функции m действительных переменных $f(x_1, x_2, \dots, x_m)$.

Функция f (*целевая функция*) минимизируется на некотором множестве $X \subset \mathbb{R}^m$. В случае, когда $X = \mathbb{R}^m$ (т.е. ограничения на переменные x_1, x_2, \dots, x_m отсутствуют) принято говорить о *задаче безусловной минимизации*. В противном случае (т.е. тогда, когда $X \neq \mathbb{R}^m$) говорят о *задаче условной минимизации*.

В данной главе рассматриваются методы решения задачи безусловной минимизации. Многие из них являются основой для перехода к более сложным методам решения задач условной минимизации.

§ 10.1. Задача безусловной минимизации функции многих переменных

1. Постановка задачи. Определения. Пусть $f(x) = f(x_1, x_2, \dots, x_m)$ — действительная функция многих переменных, определенная на множестве $X \subset \mathbb{R}^m$. Точка $\bar{x} \in X$ называется *точкой глобального минимума* функции f на множестве X , если для всех $x \in X$ выполняется неравенство $f(\bar{x}) \leq f(x)$. В этом случае значение $f(\bar{x})$ называется *минимальным значением функции* f на X .

Точка $\bar{x} \in X$ называется *точкой локального минимума* функции f , если существует такая δ -окрестность U_δ этой точки, что для всех $x \in X_\delta = X \cap U_\delta$ выполняется неравенство $f(\bar{x}) \leq f(x)$. Если же для всех $x \in X_\delta$ таких, что $x \neq \bar{x}$, выполняется строгое неравенство $f(\bar{x}) < f(x)$, то \bar{x} называется *точкой строгого локального минимума*.

Подавляющее большинство методов решения задачи безусловной минимизации в действительности являются методами поиска точки локального минимума. За исключением весьма редких случаев для нахождения точки глобального минимума, вообще говоря, не остается ничего иного, как найти все точки локального минимума и, сравнив вычис-

¹ Как и в случае одной переменной, задача максимизации сводится к задаче минимизации заменой функции f на $-f$.

ленные в этих точках значения функции f , выделить среди них точку глобального минимума. Однако такой подход связан с чрезмерно большими вычислительными затратами и вряд ли перспективен. На практике чаще используется другой подход к нахождению точки глобального минимума, который состоит в том, чтобы определить ее местоположение из анализа самой решаемой задачи, а затем применить для вычисления один из методов поиска точки локального минимума.

2. Поверхность уровня, градиент и матрица Гессе. Необходимые и достаточные условия локального минимума. Напомним некоторые определения и факты, известные из стандартного курса теории функций многих переменных.

Множество точек, для которых целевая функция принимает постоянное значение $f(\mathbf{x}) = c$, называется *поверхностью уровня*. В случае $m = 2$ это множество называют *линией уровня*. На рис. 10.1 показано, как получаются линии уровня для функции двух переменных. Функция $f(x_1, x_2)$ задает в трехмерном пространстве некоторую поверхность $u = f(x_1, x_2)$, низшая точка которой и дает решение задачи минимизации. Для того чтобы изобразить рельеф этой поверхности, проведем несколько равноотстоящих плоскостей $u = \text{const}$. Проекции на плоскость $0x_1x_2$ линий пересечения этих плоскостей с поверхностью и дают линии уровня.

Для дифференцируемой функции определен вектор из первых частных производных

$$\mathbf{g}(\mathbf{x}) = f'(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}(\mathbf{x}), \frac{\partial f}{\partial x_2}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_m}(\mathbf{x}) \right)^T,$$

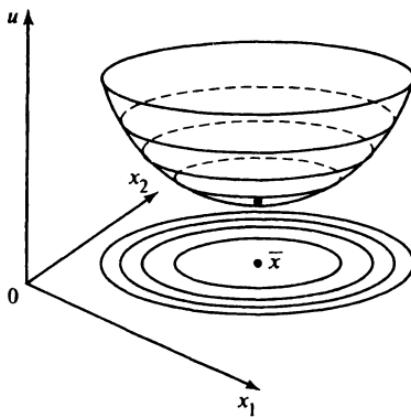


Рис. 10.1

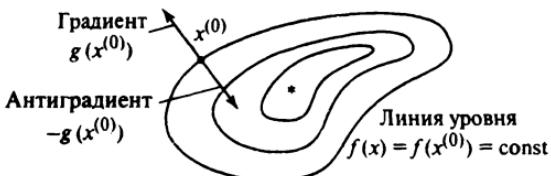


Рис. 10.2

который называется *градиентом*. Если в точке x градиент не равен нулю, то он, как известно, перпендикулярен проходящей через эту точку поверхности уровня и указывает в точке x направление наискорейшего возрастания функции. Вектор $-g(x)$ называется *антиградиентом* и указывает направление наискорейшего убывания функции (рис. 10.2).

Известно также, что равенство нулю градиента в точке x является необходимым условием того, чтобы внутренняя для множества X точка x была точкой локального минимума дифференцируемой функции f . Точка x , для которой выполняется равенство

$$f'(x) = 0, \quad (10.1)$$

называется *стационарной точкой* функции f . Равенство (10.1) представляет собой систему m нелинейных уравнений относительно компонент x_1, x_2, \dots, x_m вектора x , имеющую вид:

$$\begin{aligned} \frac{\partial f}{\partial x_1}(x_1, x_2, \dots, x_m) &= 0, \\ \dots \dots \dots \dots \dots \dots \dots & \\ \frac{\partial f}{\partial x_m}(x_1, x_2, \dots, x_m) &= 0. \end{aligned} \tag{10.2}$$

Однако не всякая стационарная точка является точкой локального минимума. Пусть функция f дважды непрерывно дифференцируема. Тогда достаточным условием того, чтобы стационарная точка x была точкой локального минимума, является положительная определенность¹ матрицы

$$G(\mathbf{x}) = f''(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_m}(\mathbf{x}) \\ \dots & \dots & \dots \\ \frac{\partial^2 f}{\partial x_m \partial x_1}(\mathbf{x}) & \dots & \frac{\partial^2 f}{\partial x_m^2}(\mathbf{x}) \end{bmatrix}, \quad (10.3)$$

¹ Определение положительно определенной симметричной матрицы см. в § 5.3.

составленной из вторых частных производных функции f . Матрицу (10.3) принято называть *матрицей Гессе*¹.

3. Выпуклые функции. Понятие выпуклости играет значительную роль в теории методов минимизации. Функция f называется *строго выпуклой*, если для любых $x \neq y$, $0 < \lambda < 1$ выполняется неравенство²

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y).$$

Это определение имеет наглядный геометрический смысл — график функции f на интервале, соединяющем точки x и y , лежит строго ниже хорды, соединяющей точки $(x, f(x))$ и $(y, f(y))$ (рис. 10.3). Для дважды непрерывно дифференцируемой функции положительная определенность матрицы Гессе $f''(x)$ является достаточным условием строгой выпуклости.

Функция f называется *сильно выпуклой* с постоянной $\kappa > 0$, если для любых x, y , $0 < \lambda < 1$ выполнено неравенство²

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \frac{\kappa}{2} \lambda(1 - \lambda)|x - y|^2. \quad (10.4)$$

Дважды непрерывно дифференцируемая функция f является сильно выпуклой тогда и только тогда, когда для всех x матрица Гессе удовлетворяет условию

$$(f''(x)\xi, \xi) \geq \kappa |\xi|^2 \quad \text{для всех } \xi \in \mathbb{R}^m,$$

где $\kappa > 0$ — постоянная, входящая в неравенство (10.4).

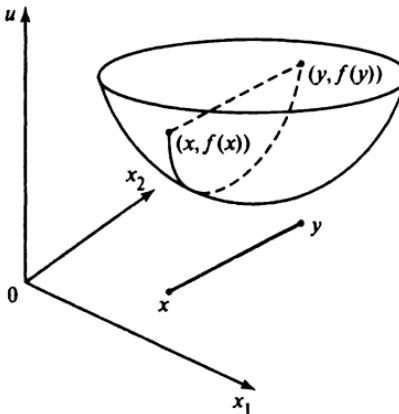


Рис. 10.3

¹ Людвиг Отто Гессе (1811—1874) — немецкий математик.

² Всюду в этой главе $|x| = \|x\|_2 = \sqrt{\sum_{i=1}^m |x_i|^2}$.

4. Задача минимизации квадратичной функции. Часто первона-
чальное исследование свойств методов безусловной минимизации про-
водится применительно к задаче минимизации квадратичной функции

$$F(x_1, x_2, \dots, x_m) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_{ij} x_i x_j - \sum_{i=1}^m b_i x_i, \quad (10.5)$$

коэффициенты a_{ij} , которой являются элементами симметричной положительно определенной матрицы A . Использовав матричные обозна-
чения, запишем функцию F так:

$$F(\mathbf{x}) = \frac{1}{2} (\mathbf{Ax}, \mathbf{x}) - (\mathbf{b}, \mathbf{x}). \quad (10.6)$$

Вычислим градиент и матрицу Гессе для функции (10.5). Дифферен-
цирование F по x_k дает

$$\frac{\partial F}{\partial x_k} = \frac{1}{2} \sum_{j=1}^m a_{kj} x_j + \frac{1}{2} \sum_{i=1}^m a_{ik} x_i - b_k.$$

Пользуясь симметрией матрицы A , получаем формулу

$$\frac{\partial F}{\partial x_k} = \sum_{j=1}^m a_{kj} x_j - b_k, \quad 1 \leq k \leq m. \quad (10.7)$$

Таким образом.

$$F'(\mathbf{x}) = \mathbf{Ax} - \mathbf{b}. \quad (10.8)$$

Дифференцируя обе части равенства (10.7) по x_l , получаем $\frac{\partial^2 F}{\partial x_l \partial x_k} = a_{lk}$. Это означает, что для квадратичной функции (10.5) матрица Гессе не зависит от \mathbf{x} и равна A .

Задача минимизации квадратичной функции представляет интерес по многим причинам. Отметим две основные из них. Во-первых, в малой окрестности точки минимума $\bar{\mathbf{x}}$ гладкая целевая функция хорошо аппроксимируется суммой первых трех слагаемых ее разло-
жения по формуле Тейлора:

$$f(\mathbf{x}) \approx F_*(\mathbf{x}) = f(\mathbf{x}^*) + (\mathbf{g}(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^*) + \frac{1}{2} (\mathbf{G}(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*), \mathbf{x} - \mathbf{x}^*) \quad (10.9)$$

с центром в точке $\mathbf{x}^* \approx \bar{\mathbf{x}}$. Функция F_* с точностью до постоянного слагаемого может быть записана в виде (10.6) с матрицей $A = \mathbf{G}(\mathbf{x}^*) \approx \mathbf{G}(\bar{\mathbf{x}})$. Поэтому можно ожидать, что вблизи точки минимума качест-
венный характер поведения последовательности $\mathbf{x}^{(n)}$, генерируемой

методом минимизации, для функции f окажется почти таким же, как и для квадратичной функции F .

Во-вторых, хорошо известно, что когда A — симметричная положительно определенная матрица, задача минимизации квадратичной функции (10.6) эквивалентна задаче решения системы линейных алгебраических уравнений

$$Ax = b. \quad (10.10)$$

Более того, решение x^0 системы (10.10) дает точку минимума функции (10.6). Действительно, $F'(x^0) = Ax^0 - b = 0$, т.е. x^0 является стационарной точкой функции F . Так как матрица Гессе $F''(x^0) = A$ положительно определена, то в точке x^0 выполнены достаточные условия минимума и, значит, $x^0 = \bar{x}$.

Таким образом, всякий метод безусловной минимизации, будучи примененным к поиску минимума функции (10.6), порождает некоторый метод решения системы (10.10).

Отметим, что поверхностями уровня квадратичной функции (10.6) служат m -мерные эллипсоиды (при $m = 2$ — эллипсы) с центром в точке \bar{x} . Отношение большей оси каждого из этих эллипсоидов (эллипсов) к меньшей оси равно $\text{cond}_2(A) = \lambda_{\max}/\lambda_{\min}$ — числу обусловленности матрицы A . Здесь λ_{\max} и λ_{\min} — максимальное и минимальное собственные значения матрицы A . Чем больше отношение $\lambda_{\max}/\lambda_{\min}$, тем сильнее вытянуты поверхности (линии) уровня.

5. Обусловленность задачи минимизации. В § 9.2 достаточно подробно обсуждалась обусловленность задачи поиска строгого локального минимума функции одной переменной. Так как ситуация в многомерном случае аналогична, то здесь изложим соответствующие положения более кратко.

Пусть \bar{x} — точка строго локального минимума дважды непрерывно дифференцируемой функции f , матрица Гессе которой в точке \bar{x} положительно определена. Предположим, что в некоторой окрестности точки минимума вычисляемые приближенные значения $f^*(x)$ удовлетворяют неравенству $|f(x) - f^*(x)| \leq \bar{\Delta} = \bar{\Delta}(f^*)$. Тогда для радиуса соответствующей области неопределенности справедлива грубая оценка $\bar{\varepsilon} \approx 2\sqrt{\bar{\Delta}/\lambda_{\min}}$, где $\lambda_{\min} > 0$ — минимальное собственное значение матрицы Гессе $f''(\bar{x})$.

Отсюда следует, что методы приближенного решения задачи минимизации, в которых используются только значения f^* , не могут, вообще говоря, гарантировать получение приближенного решения \bar{x}^* с погрешностью меньшей чем $\bar{\Delta}(\bar{x}^*) \approx 2\sqrt{\bar{\Delta}/\lambda_{\min}}$. Это означает, например, что даже для хорошо масштабированных задач (т.е. при $|\bar{x}| \sim 1$, $|f(\bar{x})| \sim 1$, $\lambda_{\min} \sim 1$) неизбежна потеря примерно половины из того числа верных значащих цифр, которые содержались в приближенных значениях $f^*(x)$.

Пусть теперь для решения задачи минимизации доступны приближенные значения градиента $(f')^*$. Тогда сведение задачи поиска точки минимума \bar{x} к эквивалентной задаче решения системы уравнений $f'(x) = 0$ существенно улучшает обусловленность задачи. В частности, справедлива такая оценка границы абсолютной погрешности:

$$\bar{\Delta}(\bar{x}^*) \approx \frac{1}{\lambda_{\min}} \bar{\Delta}((f')^*).$$

Здесь $\bar{\Delta}((f')^*)$ — граница абсолютной погрешности значений $(f')^*$, считающаяся достаточно малой величиной. Поэтому если для вычисления \bar{x} можно использовать значения градиента, то такой возможностью не следует пренебрегать.

§ 10.2. Понятие о методах спуска. Покоординатный спуск

1. Методы спуска. Большинство итерационных методов, применяемых для решения задачи безусловной минимизации функций многих переменных, относятся к классу *методов спуска*, т.е. таких методов, для которых каждая итерация (шаг) приводит к уменьшению значения целевой функции: $f(x^{(n+1)}) < f(x^{(n)})$ для всех $n \geq 0$.

Опишем структуру типичной $(n+1)$ -й итерации метода спуска в предположении, что приближение $x^{(n)}$ к точке минимума уже найдено и $x^{(n)} \neq \bar{x}$.

1°. Находят ненулевой вектор $p^{(n)}$, называемый *направлением спуска*. Этот вектор должен быть таким, чтобы при всех достаточно малых $\alpha > 0$ выполнялось неравенство

$$f(x^{(n)} + \alpha p^{(n)}) < f(x^{(n)}). \quad (10.11)$$

Если ввести функцию одной переменной α , имеющую вид

$$\varphi_n(\alpha) = f(\mathbf{x}^{(n)} + \alpha \mathbf{p}^{(n)}), \quad (10.12)$$

то неравенство (10.11) можно записать так: $\varphi_n(\alpha) < \varphi_n(0)$.

2°. Вычисляют положительное число α_n (*шаг спуска*), для которого выполняется неравенство

$$f(\mathbf{x}^{(n)} + \alpha_n \mathbf{p}^{(n)}) < f(\mathbf{x}^{(n)}), \quad (10.13)$$

или, что то же самое, неравенство $\varphi_n(\alpha_n) < \varphi_n(0)$.

3°. За очередное приближение к точке минимума принимают

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + \alpha_n \mathbf{p}^{(n)}.$$

4°. Проверяют выполнение критерия окончания итераций. Если критерий выполняется, то итерации прекращают и полагают $\bar{\mathbf{x}} \approx \mathbf{x}^{(n+1)}$. В противном случае итерации продолжают далее.

Последовательность точек $\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^{(n)}, \dots$, генерируемую методом спуска, иногда называют *траекторией спуска*.

2. Направление спуска. Заметим, что вектор $\mathbf{p}^{(n)}$ не может быть задан произвольно. В силу требования (10.11) функция $\varphi_n(\alpha)$, определяемая формулой (10.12), должна убывать в точке $\alpha = 0$. Для того чтобы это условие выполнялось, достаточно потребовать выполнения неравенства $\varphi'_n(0) < 0$. Так как $\varphi'_n(\alpha) = (f'(\mathbf{x}^{(n)} + \alpha \mathbf{p}^{(n)}), \mathbf{p}^{(n)})$, то вектор $\mathbf{p}^{(n)}$ является направлением спуска, т.е. удовлетворяет условию (10.11), если

$$(f'(\mathbf{x}^{(n)}), \mathbf{p}^{(n)}) < 0. \quad (10.14)$$

Можно показать, что для строго выпуклой функции f это неравенство выполняется тогда и только тогда, когда вектор $\mathbf{p}^{(n)}$ задает направление спуска.

Обычно именно способ выбора направления спуска $\mathbf{p}^{(n)}$ определяет конкретный численный метод, а различные варианты метода определяются далее алгоритмом вычисления шага спуска α_n . Например, выбор в качестве $\mathbf{p}^{(n)}$ антиградиента $\mathbf{p}^{(n)} = -f'(\mathbf{x}^{(n)})$ задает градиентный метод (см. § 10.3). В этом случае $(f'(\mathbf{x}^{(n)}), \mathbf{p}^{(n)}) = -|f'(\mathbf{x}^{(n)})|^2 < 0$, т.е. условие (10.14) выполняется.

3. Выбор шага спуска. Следует иметь в виду, что шаг α_n спуска определяет «истинную» величину шага $h_n = |\mathbf{x}^{(n+1)} - \mathbf{x}^{(n)}| = \alpha_n |\mathbf{p}^{(n)}|$, но совпадает с ней только тогда, когда вектор $\mathbf{p}^{(n)}$ имеет единичную длину.

В качестве α_n можно выбирать значение параметра α , которое обеспечивает достижение наименьшего значения функции f вдоль луча $x = x^{(n)} + \alpha p^{(n)}$, $\alpha \geq 0$. В этом случае для вычисления α_n требуется решение одномерной задачи минимизации функции $\varphi_n(\alpha)$. Этот путь достаточно надежен, однако может быть связан с большими вычислительными затратами.

Существуют и активно применяются другие подходы к выбору α_n , гарантирующие выполнение условия (10.13). Один из простейших подходов, называемый *дроблением шага*, состоит в следующем. Фиксируют начальное значение шага α и выбирают параметр γ , $0 < \gamma < 1$ (часто $\gamma = 1/2$). За шаг спуска принимают $\alpha_n = \alpha \gamma^{i_n}$, где i_n — первый из номеров $i = 0, 1, 2, \dots$, для которого выполнено условие

$$f(x^{(n)} + \alpha \gamma^i p^{(n)}) - f(x^{(n)}) \leq \beta \alpha \gamma^i (f'(x^{(n)}), p^{(n)}). \quad (10.15)$$

здесь $0 < \beta < 1$, β — некоторый дополнительный параметр.

4. Критерии окончания итераций. На практике часто используют следующие критерии окончания итераций:

$$|x^{(n+1)} - x^{(n)}| < \varepsilon_1, \quad (10.16)$$

$$|f(x^{(n+1)}) - f(x^{(n)})| < \varepsilon_2, \quad (10.17)$$

$$|f'(x^{(n)})| < \varepsilon_3, \quad (10.18)$$

где $\varepsilon_1, \varepsilon_2, \varepsilon_3$ — заданные положительные числа. Нередко используют различные их сочетания, например требуют, чтобы одновременно выполнялись условия (10.16) и (10.17) или даже все три условия (10.16)–(10.18).

Нередко вместо критериев (10.16), (10.17) применяют их аналоги, основанные на сочетании понятий относительной и абсолютной погрешностей:

$$|x^{(n+1)} - x^{(n)}| < \delta_1 (1 + |x^{(n+1)}|), \quad (10.19)$$

$$|f(x^{(n+1)}) - f(x^{(n)})| < \delta_2 (1 + |f(x^{(n+1)})|), \quad (10.20)$$

а также другие критерии.

К сожалению, надежные критерии окончания счета, которые были бы применимы к широкому классу задач и гарантировали бы достижение нужной точности, пока не известны.

5. Покоординатный спуск. В методе покоординатного спуска в качестве очередного направления спуска выбирают направление одной из координатных осей. Наиболее известным является *метод циклического покоординатного спуска*. Опишем очередной $(n+1)$ -й цикл одного из вариантов этого метода, в предположении, что приближение $x^{(n)}$ уже найдено.

Цикл с номером $n + 1$ состоит из m шагов. На первом шаге производят спуск по координате x_1 . Значения $x_2 = x_2^{(n)}, \dots, x_m = x_m^{(n)}$ остальных координат фиксируют, а $x_1^{(n+1)}$ выбирают из условия

$$f(x_1^{(n+1)}, x_2^{(n)}, \dots, x_m^{(n)}) = \min_{x_1} f(x_1, x_2^{(n)}, \dots, x_m^{(n)}).$$

Фактически решается задача минимизации функции одной переменной

$$f_1(x_1) = f(x_1, x_2^{(n)}, \dots, x_m^{(n)}).$$

На втором шаге производят спуск по координате x_2 . Значения $x_1 = x_1^{(n+1)}, x_3 = x_3^{(n)}, \dots, x_m = x_m^{(n)}$ остальных координат фиксируют и $x_2^{(n+1)}$ выбирают как решение задачи одномерной минимизации

$$f(x_1^{(n+1)}, x_2^{(n+1)}, x_3^{(n)}, \dots, x_m^{(n)}) = \min_{x_2} f(x_1^{(n+1)}, x_2, x_3^{(n)}, \dots, x_m^{(n)}).$$

Аналогично осуществляют остальные шаги. На последнем m -м шаге координату $x_m^{(n+1)}$ определяют из условия

$$f(x_1^{(n+1)}, \dots, x_{m-1}^{(n+1)}, x_m^{(n+1)}) = \min_{x_m} f(x_1^{(n+1)}, \dots, x_{m-1}^{(n+1)}, x_m).$$

В результате получается очередное приближение $x^{(n+1)}$ к точке минимума. Далее цикл метода снова повторяют.

На рис. 10.4 изображена геометрическая иллюстрация циклического покоординатного спуска.

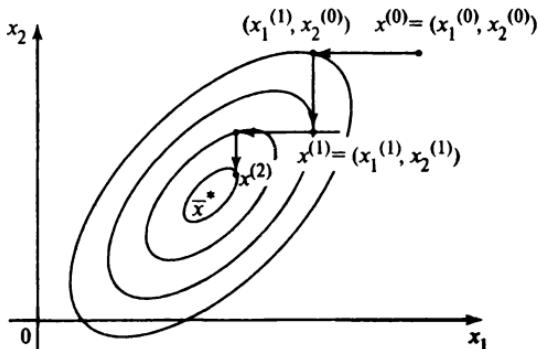


Рис. 10.4

Применим метод циклического покоординатного спуска для минимизации квадратичной функции (10.5). Так как на k -м шаге очередного цикла значение координаты $x_k^{(n+1)}$ определяется из условия минимизации функции F по направлению x_k , то необходимо, чтобы в точке $(x_1^{(n+1)}, \dots, x_{k-1}^{(n+1)}, x_k^{(n+1)}, x_{k+1}^{(n)}, \dots, x_m^{(n)})$ производная $\partial F / \partial x_k$ обращалась в нуль. Учитывая формулу (10.7), получаем уравнение

$$\sum_{j=1}^k a_{kj} x_j^{(n+1)} + \sum_{j=k+1}^m a_{kj} x_j^{(n)} - b_k = 0,$$

откуда находим

$$x_k^{(n+1)} = a_{kk}^{-1} \left(b_k - \sum_{j=1}^{k-1} a_{kj} x_j^{(n+1)} - \sum_{j=k+1}^m a_{kj} x_j^{(n)} \right). \quad (10.21)$$

Последовательные вычисления по формуле (10.21) для $k = 1, 2, \dots, m$ и составляют один цикл метода покоординатного спуска. Заметим, что этот метод применительно к решению системы линейных алгебраических уравнений $Ax = b$ с симметричной положительно определенной матрицей дает известный итерационный метод Зейделя (см. гл. 6), сходящийся со скоростью геометрической прогрессии.

§ 10.3. Градиентный метод

Рассмотрим задачу безусловной минимизации дифференцируемой функции многих переменных $f(\mathbf{x})$. Пусть $\mathbf{x}^{(n)}$ — приближение к точке минимума $\bar{\mathbf{x}}$, а $\mathbf{g}^{(n)} = g(\mathbf{x}^{(n)})$ — значение градиента в точке $\mathbf{x}^{(n)}$. Выше уже отмечалось, что в малой окрестности точки $\mathbf{x}^{(n)}$ направление наискорейшего убывания функции f задается антиградиентом $-\mathbf{g}^{(n)}$. Это свойство существенно используется в ряде методов минимизации. В рассматриваемом ниже *градиентном методе* за направление спуска из точки $\mathbf{x}^{(n)}$ непосредственно выбирается $\mathbf{p}^{(n)} = -\mathbf{g}^{(n)}$. Таким образом, согласно градиентному методу

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \alpha_n \mathbf{g}^{(n)}. \quad (10.22)$$

Существуют различные способы выбора шага α_n , каждый из которых задает определенный вариант градиентного метода.

1. Метод наискорейшего спуска. Рассмотрим функцию $\varphi_n(\alpha) = f(\mathbf{x}^{(n)} - \alpha \mathbf{g}^{(n)})$ одной скалярной переменной $\alpha \geq 0$ и выберем в качестве α_n то значение, для которого выполняется равенство

$$\varphi_n(\alpha_n) = \min_{0 < \alpha} \varphi_n(\alpha). \quad (10.23)$$

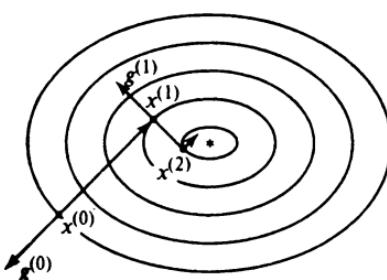


Рис. 10.5

Этот метод, предложенный в 1845 г. О. Коши¹, принято теперь называть *методом наискорейшего спуска*.

На рис. 10.5 изображена геометрическая иллюстрация этого метода для минимизации функции двух переменных. Из начальной точки $x^{(0)}$ перпендикулярно линии уровня $f(x) = f(x^{(0)})$ в направлении $p^{(0)} = -g^{(0)}$ спуск продолжают до тех пор, пока не будет достигнуто минимальное вдоль луча $x^{(0)} + \alpha p^{(0)}$ ($\alpha > 0$) значение функции f . В найденной точке $x^{(1)}$ этот луч касается линии уровня $f(x) = f(x^{(1)})$. Затем из точки $x^{(1)}$ проводят спуск в перпендикулярном линии уровня направлении $p^{(1)} = -g^{(1)}$ до тех пор, пока соответствующий луч не коснется в точке $x^{(2)}$, проходящей через эту точку линии уровня, и т.д.

Отметим, что на каждой итерации выбор шага α_n предполагает решение задачи одномерной минимизации (10.23). Иногда эту операцию удается выполнить аналитически, например для квадратичной функции.

Применим метод наискорейшего спуска для минимизации квадратичной функции

$$F(\mathbf{x}) = \frac{1}{2} (\mathbf{A}\mathbf{x}, \mathbf{x}) - (\mathbf{b}, \mathbf{x}) \quad (10.24)$$

с симметричной положительно определенной матрицей \mathbf{A} .

Согласно формуле (10.8), в этом случае $\mathbf{g}^{(n)} = \mathbf{A}\mathbf{x}^{(n)} - \mathbf{b}$. Поэтому формула (10.22) выглядит здесь так:

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \alpha_n (\mathbf{A}\mathbf{x}^{(n)} - \mathbf{b}). \quad (10.25)$$

Заметим, что

$$\begin{aligned} \varphi_n(\alpha) &= \frac{1}{2} (\mathbf{A}(\mathbf{x}^{(n)} - \alpha\mathbf{g}^{(n)}), \mathbf{x}^{(n)} - \alpha\mathbf{g}^{(n)}) - (\mathbf{b}, \mathbf{x}^{(n)} - \alpha\mathbf{g}^{(n)}) = \\ &= \frac{1}{2} (\mathbf{A}\mathbf{g}^{(n)}, \mathbf{g}^{(n)})\alpha^2 - (\mathbf{g}^{(n)}, \mathbf{g}^{(n)})\alpha + \frac{1}{2} (\mathbf{A}\mathbf{x}^{(n)}, \mathbf{x}^{(n)}) - (\mathbf{b}, \mathbf{x}^{(n)}). \end{aligned}$$

¹ Огюстен Луи Коши (1789—1857) — французский математик, один из создателей современного математического анализа, теории дифференциальных уравнений и др.

Эта функция является квадратичной функцией параметра α и достигает минимума при таком значении $\alpha = \alpha_n$, для которого

$$\varphi'_n(\alpha_n) = (\mathbf{A}\mathbf{g}^{(n)}, \mathbf{g}^{(n)})\alpha_n - (\mathbf{g}^{(n)}, \mathbf{g}^{(n)}) = 0.$$

Таким образом, применительно к минимизации квадратичной функции (10.24) метод наискорейшего спуска эквивалентен расчету по формуле (10.25), где

$$\alpha_n = \frac{(\mathbf{g}^{(n)}, \mathbf{g}^{(n)})}{(\mathbf{A}\mathbf{g}^{(n)}, \mathbf{g}^{(n)})}. \quad (10.26)$$

Замечание 1. Поскольку точка минимума функции (10.24) совпадает с решением системы $\mathbf{Ax} = \mathbf{b}$, метод наискорейшего спуска (10.25), (10.26) может применяться и как итерационный метод решения систем линейных алгебраических уравнений с симметричными положительно определенными матрицами (см. § 6.4, пункт 8). При этом $\mathbf{g}^{(n)} = -\mathbf{r}^{(n)}$, где $\mathbf{r}^{(n)} = \mathbf{b} - \mathbf{Ax}^{(n)}$ — невязка, отвечающая приближению $\mathbf{x}^{(n)}$.

Замечание 2. Отметим, что $\alpha_n^{-1} = \rho(\mathbf{g}^{(n)})$, где $\rho(\mathbf{x}) = (\mathbf{Ax}, \mathbf{x})/(\mathbf{x}, \mathbf{x})$ — отношение Рэлея (см. § 8.1).

Пример 10.1. Применим метод наискорейшего спуска для минимизации квадратичной функции $f(x_1, x_2) = x_1^2 + 2x_2^2 - 4x_1 - 4x_2$.

Заметим, что $f(x_1, x_2) = (x_1 - 2)^2 + 2(x_2 - 1)^2 - 6$. Поэтому точное значение $\bar{\mathbf{x}} = (2, 1)^T$ точки минимума нам заранее известно. Запишем данную функцию в виде (10.24), где матрица $\mathbf{A} = \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix}$ и вектор $\mathbf{b} = (4, 4)^T$. Как

нетрудно видеть, $\mathbf{A} = \mathbf{A}^T > 0$.

Возьмем начальное приближение $\mathbf{x}^{(0)} = (0, 0)^T$ и будем вести вычисления по формулам (10.25), (10.26).

1-я итерация.

$$\begin{aligned} \mathbf{g}^{(0)} &= \mathbf{Ax}^{(0)} - \mathbf{b} = (-4, -4)^T, \quad \alpha_0 = \frac{(\mathbf{g}^{(0)}, \mathbf{g}^{(0)})}{(\mathbf{A}\mathbf{g}^{(0)}, \mathbf{g}^{(0)})} = \frac{4^2 + 4^2}{2 \cdot 4^2 + 4 \cdot 4^2} = \frac{1}{3}, \\ \mathbf{x}^{(1)} &= \mathbf{x}^{(0)} - \alpha_0 \mathbf{g}^{(0)} = \left(\frac{4}{3}, \frac{4}{3} \right)^T. \end{aligned}$$

2-я итерация.

$$\mathbf{g}^{(1)} = \mathbf{Ax}^{(1)} - \mathbf{b} = \left(-\frac{4}{3}, \frac{4}{3} \right)^T, \quad \alpha_1 = \frac{(\mathbf{g}^{(1)}, \mathbf{g}^{(1)})}{(\mathbf{A}\mathbf{g}^{(1)}, \mathbf{g}^{(1)})} = \frac{(4/3)^2 + (4/3)^2}{2 \cdot (4/3)^2 + 4 \cdot (4/3)^2} = \frac{1}{3},$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \alpha_1 \mathbf{g}^{(1)} = \left(\frac{16}{9}, \frac{8}{9} \right)^T.$$

Можно показать, что для всех $n \geq 1$ на n -й итерации будут получены значения

$$\mathbf{g}^{(n)} = \frac{-4}{3^n} (1, (-1)^n)^T, \quad \alpha_n = \frac{1}{3}, \quad \mathbf{x}^{(n)} = \bar{\mathbf{x}} - \frac{1}{3^n} (2, (-1)^n)^T.$$

Заметим, что $|\mathbf{x}^{(n)} - \bar{\mathbf{x}}| = \frac{\sqrt{5}}{3^n} \rightarrow 0$ при $n \rightarrow \infty$. Таким образом, последовательность $\mathbf{x}^{(n)}$, полученная методом наискорейшего спуска, сходится со скоростью геометрической прогрессии, знаменатель которой $q = 1/3$.

На рис. 10.5 изображена именно траектория спуска, которая была получена в данном примере. ▲

Для случая минимизации квадратичной функции справедлив следующий общий результат [20].

Теорема 10.1. Пусть \mathbf{A} — симметричная положительно определенная матрица и минимизируется квадратичная функция (10.24). Тогда при любом выборе начального приближения метод наискорейшего спуска (10.25), (10.26) сходится и верна следующая оценка погрешности:

$$|\mathbf{x}^{(n)} - \mathbf{x}| \leq \sqrt{\text{cond}_2(\mathbf{A})} q^n |\mathbf{x}^{(0)} - \mathbf{x}| \quad (10.27)$$

Здесь $q = \frac{\text{cond}_2(\mathbf{A}) - 1}{\text{cond}_2(\mathbf{A}) + 1} < 1$, где $\text{cond}_2(\mathbf{A}) = \frac{\lambda_{\max}}{\lambda_{\min}}$, а λ_{\min} и λ_{\max} — минимальное и максимальное собственные значения матрицы \mathbf{A} . ■

Отметим что этот метод сходится со скоростью геометрической прогрессии, знаменатель которой равен q . Если λ_{\min} и λ_{\max} близки, то $\text{cond}_2(\mathbf{A}) \sim 1$ и q мало. В этой ситуации метод сходится достаточно быстро. Например, в примере 10.1 имеем $\lambda_{\min} = 2$, $\lambda_{\max} = 4$ и поэтому $\text{cond}_2(\mathbf{A}) = 2$, $q = (2 - 1)/(2 + 1) = 1/3$. Если же $\lambda_{\min} \ll \lambda_{\max}$, то $\text{cond}_2(\mathbf{A}) \gg 1$, $q \approx 1$ и следует ожидать медленной сходимости метода наискорейшего спуска.

Пример 10.2. Применение метода наискорейшего спуска для минимизации квадратичной функции $f(x_1, x_2) = x_1^2 + 10x_2^2 - 4x_1 - 4x_2$ при начальном приближении $\mathbf{x}^{(0)} = (0, 0)^T$ дает последовательность приближений $\mathbf{x}^{(n)} = \bar{\mathbf{x}} - 2(9/11)^n (1, (-1)^n)^T$, где $\bar{\mathbf{x}} = (2, 0.2)^T$. Траектория спуска изображена на рис. 10.6.

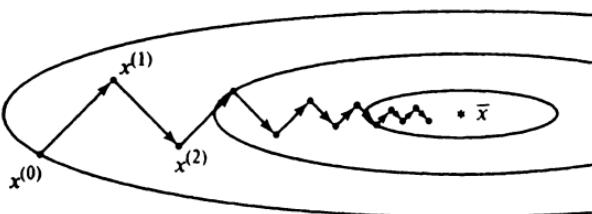


Рис. 10.6

Последовательность $x^{(n)}$ сходится здесь со скоростью геометрической прогрессии, знаменатель которой равен $q = 9/11$, т.е. существенно медленнее, чем в примере 10.1. Так как здесь $A = \begin{bmatrix} 2 & 0 \\ 0 & 20 \end{bmatrix}$, то $\lambda_{\min} = 2$, $\lambda_{\max} = 20$, $\text{cond}_2(A) = 20/2 = 10$, $q = (10 - 1)/(10 + 1) = 9/11$, и полученный результат вполне согласуется с оценкой (10.27). \blacktriangle

Замечание 1 Мы сформулировали теорему о сходимости метода наискорейшего спуска в предположении, что целевая функция является квадратичной. В общем случае, если минимизируемая функция строго выпуклая и имеет точку минимума \bar{x} , то также независимо от выбора начального приближения полученная указанным методом последовательность $x^{(n)}$ сходится к \bar{x} при $n \rightarrow \infty$. При этом после попадания $x^{(n)}$ в достаточно малую окрестность точки минимума сходимость становится линейной и знаменатель соответствующей геометрической прогрессии оценивается сверху величиной $\bar{q} \approx \frac{x-1}{x+1}$, где $x = \frac{\lambda_{\max}}{\lambda_{\min}}$ — отношение максимального и минимального собственных значений матрицы Гессе $G(\bar{x})$, т.е. $x = \text{cond}_2(G(\bar{x}))$.

Замечание 2. Для квадратичной целевой функции (10.24) решение задачи одномерной минимизации (10.23) удается найти в виде простой явной формулы (10.26). Однако для большинства других нелинейных функций этого сделать нельзя и для вычисления a_n методом наискорейшего спуска приходится применять численные методы одномерной минимизации типа тех, которые были рассмотрены в гл. 9.

2. Проблема «оврагов». Из проведенного выше обсуждения следует, что градиентный метод сходится достаточно быстро, если для минимизируемой функции поверхности уровня близки к сферам (при $m = 2$ линии

уровня близки к окружностям). Для таких функций $x = \lambda_{\max}/\lambda_{\min} \approx 1$. Теорема 10.1, замечание 1, а также результат примера 10.2 указывают на то, что скорость сходимости резко падает при увеличении величины x . Действительно, известно, что градиентный метод сходится очень медленно, если поверхности уровня минимизируемой функции сильно вытянуты в некоторых направлениях. В двумерном случае рельеф соответствующей поверхности $u = f(x_1, x_2)$ напоминает рельеф местности с оврагом (рис. 10.7). Поэтому такие функции принято называть *овражными*. Вдоль направлений, характеризующих «дно оврага», овражная функция меняется незначительно, а в других направлениях, характеризующих «склон оврага», происходит резкое изменение функции.

Если начальная точка $x^{(0)}$ попадает на «склон оврага», то направление градиентного спуска оказывается почти перпендикулярным «дну оврага» и очередное приближение $x^{(1)}$ попадает на противоположный «склон оврага». Следующий шаг в направлении ко «дну оврага» возвращает приближение $x^{(2)}$ на первоначальный «склон оврага». В результате вместо того чтобы двигаться вдоль «дна оврага» в направлении к точке минимума, траектория спуска совершает зигзагообразные скачки поперек «оврага», почти не приближаясь к цели (рис. 10.7).

Для ускорения сходимости градиентного метода при минимизации овражных функций разработан ряд специальных «овражных» методов. Дадим представление об одном из простейших приемов. Из двух близких начальных точек $z^{(0)}$ и $z^{(1)}$ совершают градиентный спуск на «дно оврага». Через найденные точки $x^{(0)}$ и $x^{(1)}$ проводят прямую, вдоль которой совершают большой «овражный» шаг (рис. 10.8). Из найденной таким образом точки $z^{(2)}$ снова делают один шаг градиентного спуска в точку $x^{(2)}$. Затем совершают второй «овражный» шаг вдоль прямой, проходящей через точки $x^{(1)}$ и $x^{(2)}$, и т.д. В результате движение вдоль «дна оврага» к точке минимума существенно ускоряется.

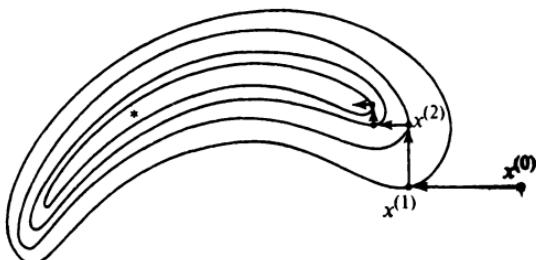


Рис. 10.7

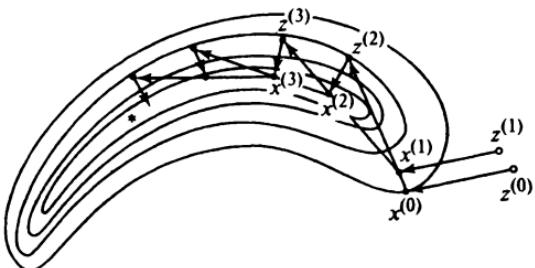


Рис. 10.8

Более подробную информацию о проблеме «оврагов» и «овражных» методах можно найти, например в [9, 20].

3. Другие подходы к определению шага спуска. Как нетрудно понять, на каждой итерации было бы желательно выбирать направление спуска $p^{(n)}$, близкое к тому направлению, перемещение вдоль которого приводит из точки $x^{(n)}$ в точку \bar{x} . К сожалению, антиградиент $p^{(n)} = -g^{(n)}$ является, как правило, неудачным направлением спуска. Особенно ярко это проявляется для овражных функций. Поэтому возникает сомнение в целесообразности тщательного поиска решения задачи одномерной минимизации (10.23) и появляется желание сделать в направлении $p^{(n)}$ лишь такой шаг, который бы обеспечил «существенное убывание» функции f . Более того, на практике иногда довольноствуются определением значения $\alpha_n > 0$, которое просто обеспечивает уменьшение значения целевой функции.

В одном из простейших алгоритмов (типа дробления шага) такого выбора шага α_n фиксируют начальное значение $\alpha > 0$ и значение параметра γ , $0 < \gamma < 1$. За α_n принимают $\alpha_n = \alpha \cdot \gamma^{i_n}$, где i_n — первый из номеров $i = 0, 1, 2, \dots$, для которого выполнено условие убывания

$$f(x^{(n)} - \alpha \gamma^i g^{(n)}) - f(x^{(n)}) < 0. \quad (10.28)$$

Однако при таком выборе α_n нет гарантии, что последовательность $x^{(n)}$ будет сходиться к точке минимума даже для простой квадратичной функции (10.24). Условие (10.28) является слишком слабым: последовательность $x^{(n)}$, незначительно уменьшая значения функции f , может «останавливаться», не доходя до точки \bar{x} . Такое поведение последова-

тельности $\mathbf{x}^{(n)}$ можно предотвратить, если заменить условие (10.28) условием «существенного убывания» функции f на каждой итерации:

$$f(\mathbf{x}^{(n)} - \alpha\gamma'\mathbf{g}^{(n)}) - f(\mathbf{x}^{(n)}) \leq -\beta\alpha\gamma'|\mathbf{g}^{(n)}|^2. \quad (10.29)$$

Здесь β ($0 < \beta < 1$) — дополнительный параметр.

Заметим, что для рассматриваемого градиентного метода $\mathbf{p}^{(n)} = -\mathbf{g}^{(n)} = -f'(\mathbf{x}^{(n)})$ и поэтому неравенство (10.29) в точности совпадает с неравенством (10.15), используемым в методах спуска при дроблении шага.

Пример 10.3. Продемонстрируем применение градиентного метода с дроблением шага к задаче минимизации квадратичной функции $f(x_1, x_2) = x_1^2 + 2x_2^2 - 4x_1 - 4x_2$ из примера 10.1. Для выбора значения шага будем использовать условие (10.29). Воспользуемся следующими краткими обозначениями: $\alpha_i = \alpha\gamma^i$, $\mathbf{x}^{(n, i)} = \mathbf{x}^{(n)} - \alpha_i\mathbf{g}^{(n)}$.

Заметим, что $\mathbf{g}^{(n)} = (2x_1^{(n)} - 4, 4x_2^{(n)} - 4)^T$.

Выберем начальное приближение $\mathbf{x}^{(n)} = (0, 0)^T$, начальное значение шага $\alpha = \alpha_0 = 1$, значения параметров $\gamma = 1/2$, $\beta = 3/4$. Вычислим значения $f(\mathbf{x}^{(0)}) = 0$, $\mathbf{g}^{(0)} = (-4, -4)^T$.

1-я итерация. Вычисляем $\mathbf{x}^{(0, 0)} = \mathbf{x}^{(0)} - \alpha_0\mathbf{g}^{(0)} = (4, 4)^T$, $f(\mathbf{x}^{(0, 0)}) = 16$. Так как значение функции не уменьшилось, то следует уменьшить шаг: $\alpha_1 = \alpha_0/2 = 0.5$.

Вычисляем $\mathbf{x}^{(0, 1)} = \mathbf{x}^{(0)} - \alpha_1\mathbf{g}^{(0)} = (2, 2)^T$, $f(\mathbf{x}^{(0, 1)}) = -4$. Поскольку $f(\mathbf{x}^{(0, 1)}) - f(\mathbf{x}^{(0)}) = -4 > -\beta\alpha_1|\mathbf{g}^{(0)}|^2 = -12$, условие (10.29) не выполняется и следует снова уменьшить шаг: $\alpha_2 = \alpha_1/2 = 0.25$.

Вычисляем $\mathbf{x}^{(0, 2)} = \mathbf{x}^{(0)} - \alpha_2\mathbf{g}^{(0)} = (1, 1)^T$, $f(\mathbf{x}^{(0, 2)}) = -5$. Имеем $f(\mathbf{x}^{(0, 2)}) - f(\mathbf{x}^{(0)}) = -5 > -\beta\alpha_2|\mathbf{g}^{(0)}|^2 = -6$, т.е. условие (10.29) не выполняется. Уменьшаем шаг: $\alpha_3 = \alpha_2/2 = 0.125$.

Вычисляем $\mathbf{x}^{(0, 3)} = \mathbf{x}^{(0)} - \alpha_3\mathbf{g}^{(0)} = (0.5, 0.5)^T$, $f(\mathbf{x}^{(0, 3)}) = -3.25$. Так как $f(\mathbf{x}^{(0, 3)}) - f(\mathbf{x}^{(0)}) = -3.25 < -\beta\alpha_3|\mathbf{g}^{(0)}|^2 = -3$, то условие (10.29) выполнено.

Положим $\mathbf{x}^{(1)} = (0.5, 0.5)^T$; напомним, что $f(\mathbf{x}^{(1)}) = -3.25$. Вычислим $\mathbf{g}^{(1)} = (-3, -2)^T$ и положим $\alpha_0 = 1$.

Далее вычисления следует продолжить до выполнения какого-либо принятого критерия окончания итераций. ▲

4. Влияние погрешности вычислений. Один из существенных недостатков градиентного метода связан с его чувствительностью к погрешностям вычислений. Особенно сильно этот недостаток сказы-

вается в малой окрестности точки минимума, где антиградиент, задающий направление поиска, мал по модулю. Поэтому эффективность градиентного метода на завершающей стадии поиска существенно ниже, чем на начальной стадии.

§ 10.4. Метод Ньютона

1. Простейший вариант метода Ньютона и метод Ньютона с дроблением шага. Пусть в некоторой окрестности точки минимума \bar{x} функция f является сильно выпуклой и дважды непрерывно дифференцируемой. Далее, пусть $x^{(n)}$ — хорошее приближение к \bar{x} . Тогда в малой окрестности точки $x^{(n)}$ функция f достаточно точно аппроксимируется квадратичной функцией

$$F_n(x) = f(x^{(n)}) + (\mathbf{g}^{(n)}, x - x^{(n)}) + \frac{1}{2} (\mathbf{G}^{(n)}(x - x^{(n)}), x - x^{(n)}),$$

являющейся суммой первых трех членов ее разложения по формуле Тейлора (ср. с формулой (10.9)). Здесь $\mathbf{g}^{(n)} = f'(x^{(n)})$, $\mathbf{G}^{(n)} = f''(x^{(n)})$.

Можно ожидать, что точка $\bar{x}^{(n)}$, в которой достигается минимум функции F_n , будет значительно лучшим приближением к \bar{x} , чем $x^{(n)}$. Учитывая, что $F'_n(x) = \mathbf{g}^{(n)} + \mathbf{G}^{(n)}(x - x^{(n)})$, замечаем, что точка $\bar{x}^{(n)}$ может быть определена из необходимого условия экстремума:

$$\mathbf{g}^{(n)} + \mathbf{G}^{(n)}(\bar{x}^{(n)} - x^{(n)}) = 0.$$

Таким образом, чтобы попасть из точки $x^{(n)}$ в точку $\bar{x}^{(n)}$, нужно переместиться вдоль вектора $\mathbf{p}^{(n)} = \bar{x}^{(n)} - x^{(n)}$, который определяется из системы линейных алгебраических уравнений

$$\mathbf{G}^{(n)} \mathbf{p}^{(n)} = -\mathbf{g}^{(n)}. \quad (10.30)$$

Вектор $\mathbf{p}^{(n)}$ принято называть *ニュтоновским направлением*, а метод спуска

$$x^{(n+1)} = x^{(n)} + \alpha_n \mathbf{p}^{(n)} \quad (10.31)$$

с таким выбором $\mathbf{p}^{(n)}$ — *методом Ньютона*.

Отметим, что ньютоновское направление является направлением спуска. В самом деле, в силу равенства (10.30) для $\mathbf{p}^{(n)}$ верна формула $\mathbf{p}^{(n)} = -[\mathbf{G}^{(n)}]^{-1} \mathbf{g}^{(n)}$. Матрица $[\mathbf{G}^{(n)}]^{-1}$ положительно определена (это следует из положительной определенности матрицы $\mathbf{G}^{(n)}$). Поэтому

$$(f'(x^{(n)}), \mathbf{p}^{(n)}) = -([\mathbf{G}^{(n)}]^{-1} \mathbf{g}^{(n)}, \mathbf{g}^{(n)}) < 0.$$

Таким образом, условие (10.14) выполняется и $p^{(n)}$ действительно задает направление спуска.

Различные варианты метода (10.31), (10.30) связаны с различными способами выбора шагов α_n . Заметим, что при выборе $\alpha_n = 1$ рассматриваемый метод в точности совпадает с методом Ньютона решения систем нелинейных уравнений, примененным к решению системы $f'(x) = 0$. Отсюда — и название метода.

Простым следствием теоремы 7.3 о сходимости метода Ньютона решения систем нелинейных уравнений является следующая теорема.

Теорема 10.2. Пусть в некоторой окрестности U точки минимума \bar{x} функция f является сильно выпуклой и трижды непрерывно дифференцируемой. Тогда найдется такая малая δ -окрестность точки \bar{x} , что при произвольном выборе начального приближения $x^{(0)}$ из этой окрестности последовательность $x^{(n)}$, вычисляемая с помощью метода (10.30), (10.31) при $\alpha_n = 1$, не выходит за пределы δ -окрестности точки \bar{x} и сходится к ней квадратично. ■

Замечание. Квадратичная скорость сходимости метода позволяет использовать простой практический критерий окончания:

$$|x^{(n+1)} - x^{(n)}| < \varepsilon. \quad (10.32)$$

Теорема 10.2 указывает на то, что метод Ньютона сходится очень быстро, и практика вычислений это подтверждает. Однако существенным недостатком рассмотренного варианта метода является необходимость выбора достаточно хорошего начального приближения, которое на начальной стадии поиска точки минимума, как правило, отсутствует. Поэтому метод Ньютона с выбором $\alpha_n = 1$ чаще применяют на завершающем этапе поиска \bar{x} , когда с помощью других методов уже найдено достаточно точное приближение к точке минимума.

Указанного недостатка в значительной степени лишен вариант метода Ньютона, в котором в качестве шага спуска выбирается $\alpha_n = \gamma^{i_n}$, где i_n — первый среди номеров $i \geq 0$, для которых выполняется неравенство

$$f(x^{(n)} + \gamma^i p^{(n)}) - f(x^{(n)}) \leq \beta \gamma^i (g^{(n)}, p^{(n)}).$$

Здесь $0 < \gamma < 1$, $0 < \beta < 1/2$ — параметры метода. Метод Ньютона с таким выбором α_n для широкого класса функций сходится при любом выборе начального приближения $x^{(0)} \in U$ и этим выгодно отличается от метода с выбором $\alpha_n = 1$.

Теорема 10.3. Пусть трижды непрерывно дифференцируемая в \mathbb{R}^m функция f имеет точку минимума \bar{x} и ее матрица Гессе $f''(x)$ положительно определена. Тогда при любом начальном приближении $x^{(0)}$ последовательность $x^{(n)}$, вычисляемая методом Ньютона с выбором $\alpha_n = \gamma^n$, сходится к \bar{x} с квадратичной скоростью. Более того, найдется номер n_0 такой, что для всех $n \geq n_0$ выполняется равенство $\alpha_n = 1$. ■

Замечание Используют и другие способы выбора α_n . Например, иногда α_n выбирают из условия $\varphi_n(\alpha_n) = \min_{0 \leq \alpha \leq 1} \varphi_n(\alpha)$, где

$$\varphi_n(\alpha) = f(x^{(n)} + \alpha p^{(n)}).$$

Квадратичная скорость сходимости метода Ньютона, а также возможность использования матрицы Гессе для контроля за соблюдением достаточных условий экстремума делают этот метод чрезвычайно привлекательным при решении задачи безусловной минимизации.

Пример 10.4. Применим метод Ньютона (10.30), (10.31) с $\alpha_n = 1$ для поиска точки минимума функции $f(x_1, x_2) = x_1^2 + 10(x_2 - \sin x_1)^2$ с точностью $\epsilon = 10^{-5}$. Будем использовать критерий окончания итераций (10.32) и вести вычисления на 6-разрядном десятичном компьютере. Отметим, что минимальное и равное нулю значение функция f достигает в точке $\bar{x} = (0, 0)^T$.

Имеем:

$$\begin{aligned} \mathbf{g}^{(n)} &= f'(x_1^{(n)}, x_2^{(n)}) = \\ &= (2x_1^{(n)} + 20(\sin x_1^{(n)} - x_2^{(n)}) \cdot \cos x_1^{(n)}, 20(x_2^{(n)} - \sin x_1^{(n)}))^T. \end{aligned}$$

$$\mathbf{G}^{(n)} = f''(x^{(n)}) = \begin{bmatrix} 2 + 20(\cos 2x_1^{(n)} + x_2^{(n)} \cdot \sin x_1^{(n)}) & -20 \cos x_1^{(n)} \\ -20 \cos x_1^{(n)} & 20 \end{bmatrix}.$$

Возьмем за начальное приближение $x^{(0)} = (1, 1)^T$

1-я итерация. Вычислим

$$\mathbf{g}^{(0)} = \begin{bmatrix} 0.286928 \\ 3.17058 \end{bmatrix}, \quad \mathbf{G}^{(0)} = \begin{bmatrix} 10.5065 & -10.8060 \\ -10.8060 & 20 \end{bmatrix}.$$

Таблица 10.1

n	$x_1^{(n)}$	$x_2^{(n)}$	$ x^{(n)} - x^{(n-1)} $
0	1.00000000	1.00000000	
1	0.57155400	0.60997200	$6 \cdot 10^{-1}$
2	0.15541900	0.19094400	$6 \cdot 10^{-1}$
3	0.00824200	0.00939200	$2 \cdot 10^{-1}$
4	0.00000082	0.00000101	$1 \cdot 10^{-2}$
5	0.00000000	0.00000922	$2 \cdot 10^{-6}$

Таким образом, при $n = 0$ система (10.30) принимает следующий вид:

$$\begin{aligned} 10.5065p_1^{(0)} - 10.8060p_2^{(0)} &= -0.286928, \\ -10.8060p_1^{(0)} + 20p_2^{(0)} &= -3.17058. \end{aligned}$$

Решая ее, получаем $p_1^{(0)} = -0.428445$, $p_2^{(0)} = -0.390018$. Далее в соответствии с формулой (10.31) полагаем $x_1^{(1)} = x_1^{(0)} + p_1^{(0)} = 0.571555$, $x_2^{(1)} = x_2^{(0)} + p_2^{(0)} = 0.609982$. Так как $|x^{(1)} - x^{(0)}| \approx 0.6 > \varepsilon$, то вычисления следует продолжить далее.

Результаты следующих итераций приведены в табл. 10.1.

При $n = 5$ критерий окончания выполняется и вычисления следует прекратить. Итак, $\bar{x} \approx (0.00000, 0.00000)^T$. ▲

2. Понятие о квазиньютоновских методах. Высокая скорость сходимости метода Ньютона связана с тем, что в нем используется матрица Гессе, содержащая информацию о кривизнах функции f . В результате минимизируемую функцию удается достаточно точно аппроксимировать последовательностью квадратичных функций F_n . В методах, которые называют *квазиньютоновскими*¹, также фактически используется квадратичная аппроксимация функции f , но в отличие от метода Ньютона в них не применяется вычисление вторых производных. Информация о кривизнах функции f здесь накапливается постепенно и является результатом наблюдений за изменением градиента $g^{(n)}$.

¹ Первоначально эти методы называли методами переменной метрики.

Направление спуска в квазиньютоновских методах определяется как решение системы уравнений

$$\mathbf{B}^{(n)} \mathbf{p}^{(n)} = -\mathbf{g}^{(n)},$$

в которой $\mathbf{B}^{(n)}$ — текущее приближение к матрице Гессе. Начальное приближение $\mathbf{B}^{(0)}$ обычно полагают равным единичной матрице. В таком случае $\mathbf{p}^{(0)} = -\mathbf{g}^{(0)}$ и первая итерация совпадает с одним шагом градиентного метода.

После того, как найдено приближение $\mathbf{x}^{(n+1)}$, вычисляют матрицу $\mathbf{B}^{(n+1)} = \mathbf{B}^{(n)} + \Delta\mathbf{B}^{(n)}$, где $\Delta\mathbf{B}^{(n)}$ — некоторая поправочная матрица. Во всех правилах пересчета неизменным является выполнение *квазиньютоновского условия*

$$\mathbf{B}^{(n+1)}(\mathbf{x}^{(n+1)} - \mathbf{x}^{(n)}) = \mathbf{g}^{(n+1)} - \mathbf{g}^{(n)}.$$

Известно, что при определенных условиях квазиньютоновские методы сходятся сверхлинейно. Для квадратичных же функций они дают точное значение точки минимума после конечного числа итераций, которое не превышает¹ m .

Первый из квазиньютоновских методов был предложен в 1959 г. Дэвидоном. С тех пор эти методы непрерывно совершенствовались и к настоящему времени стали одними из наиболее популярных и широко применяемых на практике. Весьма интересное и содержательное обсуждение квазиньютоновских методов содержится в книгах [28, 40, 112].

Замечание. В некоторых вариантах квазиньютоновских методов направление спуска вычисляется по формуле $\mathbf{p}^{(n)} = -\mathbf{H}^{(n)}\mathbf{g}^{(n)}$, где $\mathbf{H}^{(n)}$ — матрица, аппроксимирующая матрицу, обратную матрице Гессе.

3. Метод Ньютона и квазиньютоновские методы при наличии помех. Вычислительные ошибки (помехи) оказывают существенное влияние на поведение метода Ньютона. Природа этих ошибок может быть различной (вычисление $\mathbf{g}^{(n)}$, решение системы уравнений $\mathbf{G}^{(n)}\mathbf{p}^{(n)} = -\mathbf{g}^{(n)}$ и др.). В результате вместо ньютоновского направления $\mathbf{p}^{(n)}$ получается направление $\mathbf{p}_*^{(n)}$, вычисленное с погрешностью $\epsilon = |p^{(n)} - p_*^{(n)}|$. Как известно, метод Ньютона сходится, вообще говоря, лишь в малой окрестности U точки минимума. В случае, когда радиус этой окрестности меньше ϵ , очередное приближение $\mathbf{x}^{(n+1)} =$

¹ Естественно, что это верно только при отсутствии вычислительной погрешности.

$= \mathbf{x}^{(n)} + \alpha_n \mathbf{p}_*^{(n)}$ скорее всего окажется вне окрестности U и метод не будет сходиться. Таким образом, метод Ньютона сохраняет свои преимущества лишь при высокой точности вычислений.

Квазиньютоновские методы оказываются весьма чувствительны к ошибкам в вычислении градиента. Причина здесь состоит в том, что в основе этих методов лежит идея восстановления матрицы Гессе по результатам вычисления градиента $f'(\mathbf{x})$ в точках $\mathbf{x}^{(n)}$. Как и в методе Ньютона, здесь необходимо очень аккуратное вычисление градиента.

§ 10.5. Метод сопряженных градиентов

Метод Ньютона и квазиньютоновские методы, обсуждавшиеся в § 10.4, весьма эффективны как средство решения задач безусловной минимизации. Однако они предъявляют довольно высокие требования к объему используемой памяти компьютера. Это связано с тем, что выбор направления поиска $\mathbf{p}^{(n)}$ требует решения систем линейных алгебраических уравнений, а также с возникающей необходимостью хранения матриц типа $\mathbf{B}^{(n)}$, $\mathbf{H}^{(n)}$. Поэтому при больших m использование этих методов может оказаться невозможным. В существенной степени от этого недостатка избавлены методы сопряженных направлений.

1. Понятие о методах сопряженных направлений. Рассмотрим задачу минимизации квадратичной функции

$$F(\mathbf{x}) = \frac{1}{2} (\mathbf{A}\mathbf{x}, \mathbf{x}) - (\mathbf{b}, \mathbf{x}) \quad (10.33)$$

с симметричной положительно определенной матрицей \mathbf{A} . Напомним, что для ее решения требуется один шаг метода Ньютона и не более чем m шагов квазиньютоновского метода. Методы сопряженных направлений также позволяют найти точку минимума функции (10.33) не более чем за m шагов. Добиться этого удается благодаря специальному выбору направлений поиска.

Будем говорить, что ненулевые векторы $\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(m-1)}$ являются *взаимно сопряженными* (относительно матрицы \mathbf{A}), если $(\mathbf{A}\mathbf{p}^{(n)}, \mathbf{p}^{(l)}) = 0$ для всех $n \neq l$.

Под *методом сопряженных направлений* для минимизации квадратичной функции (10.33) будем понимать метод

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + \alpha_n \mathbf{p}^{(n)} \quad (n = 0, 1, 2, \dots, m-1), \quad (10.34)$$

в котором направления $\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(m-1)}$ взаимно сопряжены, а шаги

$$\alpha_n = -\frac{(\mathbf{g}^{(n)}, \mathbf{p}^{(n)})}{(\mathbf{A}\mathbf{p}^{(n)}, \mathbf{p}^{(n)})}, \quad \mathbf{g}^{(n)} = \mathbf{F}'(\mathbf{x}^{(n)}) = \mathbf{A}\mathbf{x}^{(n)} - \mathbf{b} \quad (10.35)$$

получаются как решение задач одномерной минимизации

$$\varphi_n(\alpha_n) = \min_{\alpha \geq 0} \varphi_n(\alpha), \quad \varphi_n(\alpha) = F(\mathbf{x}^{(n)} + \alpha \mathbf{p}^{(n)}).$$

Теорема 10.4. *Метод сопряженных направлений позволяет найти точку минимума квадратичной функции (10.33) не более чем за m шагов.* ■

Методы сопряженных направлений отличаются один от другого способом построения сопряженных направлений. Наиболее известным среди них является *метод сопряженных градиентов*.

2. Метод сопряженных градиентов. В этом методе последовательность приближений вычисляют по формулам (10.34), (10.35), а направления $\mathbf{p}^{(n)}$ строят по правилу

$$\mathbf{p}^{(0)} = -\mathbf{g}^{(0)}; \quad \mathbf{p}^{(n+1)} = -\mathbf{g}^{(n+1)} + \beta_n \mathbf{p}^{(n)}, \quad n \geq 0,$$

где

$$\beta_n = \frac{(\mathbf{A}\mathbf{p}^{(n)}, \mathbf{g}^{(n+1)})}{(\mathbf{A}\mathbf{p}^{(n)}, \mathbf{p}^{(n)})}.$$

Так как $\mathbf{p}^{(0)} = -\mathbf{g}^{(0)}$, то первый шаг этого метода совпадает с шагом метода наискорейшего спуска. Можно показать (мы этого делать не будем), что направления (10.36) действительно являются сопряженными относительно матрицы \mathbf{A} . Более того, градиенты $\mathbf{g}^{(n)}$ и $\mathbf{g}^{(k)}$ оказываются ортогональными при $n \neq k$ и взаимно сопряженными при $|n - k| > 1$. Последнему обстоятельству метод и обязан своим названием.

Пример 10.5. Применим метод сопряженных градиентов для минимизации квадратичной функции $f(x_1, x_2) = x_1^2 + 2x_2^2 - 4x_1 - 4x_2$ из примера 10.1.

Запишем f в виде $\frac{1}{2} (\mathbf{A}\mathbf{x}, \mathbf{x}) - (\mathbf{b}, \mathbf{x})$, где

$$\mathbf{A} = \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}.$$

Возьмем начальное приближение $\mathbf{x}^{(0)} = (0, 0)^T$.

1-й шаг метода совпадает с первым шагом метода наискорейшего спуска. Поэтому (см. пример 10.1) $\mathbf{g}^{(0)} = (-4, -4)^T$, $\mathbf{p}^{(0)} = -\mathbf{g}^{(0)}$, $\mathbf{x}^{(1)} = (4/3, 4/3)^T$.

2-й шаг. Вычислим $\mathbf{g}^{(1)} = \mathbf{A}\mathbf{x}^{(1)} - \mathbf{b} = (-4/3, 4/3)^T$,

$$\beta_0 = \frac{(\mathbf{A}\mathbf{p}^{(0)}, \mathbf{g}^{(0)})}{(\mathbf{A}\mathbf{p}^{(0)}, \mathbf{p}^{(0)})} = \frac{8 \cdot (-4/3) + 16 \cdot (4/3)}{8 \cdot 4 + 16 \cdot 4} = \frac{1}{9};$$

$$\mathbf{p}^{(1)} = -\mathbf{g}^{(1)} + \beta_0 \mathbf{p}^{(0)} = (4/3, -4/3)^T + (1/9)(4, 4)^T = (16/9, -8/9)^T;$$

$$\alpha_1 = -\frac{(\mathbf{g}^{(1)}, \mathbf{p}^{(1)})}{(\mathbf{A}\mathbf{p}^{(1)}, \mathbf{p}^{(1)})} = -\frac{(-4/3) \cdot (16/9) + (4/3) \cdot (-8/9)}{(32/9) \cdot (16/9) + (-32/9) \cdot (-8/9)} = \frac{3}{8};$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{p}^{(1)} = (4/3, 4/3)^T + (3/8)(16/9, -8/9)^T = (2, 1)^T.$$

Так как $\mathbf{g}^{(2)} = \mathbf{A}\mathbf{x}^{(2)} - \mathbf{b} = 0$, то $\bar{\mathbf{x}} = \mathbf{x}^{(2)} = (2, 1)^T$ и решение оказалось найденным за два шага. ▲

Для более рациональной организации вычислений расчетные формулы метода сопряженных градиентов часто переписывают в эквивалентном виде

$$\mathbf{g}^{(0)} = \mathbf{A}\mathbf{x}^{(0)} - \mathbf{b}, \quad \mathbf{p}^{(0)} = -\mathbf{g}^{(0)}, \quad (10.36)$$

$$\alpha_n = \frac{(\mathbf{g}^{(n)}, \mathbf{g}^{(n)})}{(\mathbf{A}\mathbf{p}^{(n)}, \mathbf{p}^{(n)})}, \quad (10.37)$$

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + \alpha_n \mathbf{p}^{(n)}, \quad (10.38)$$

$$\mathbf{g}^{(n+1)} = \mathbf{g}^{(n)} + \alpha_n \mathbf{A}\mathbf{p}^{(n)}, \quad (10.39)$$

$$\beta_n = \frac{(\mathbf{g}^{(n+1)}, \mathbf{g}^{(n+1)})}{(\mathbf{g}^{(n)}, \mathbf{g}^{(n)})}, \quad (10.40)$$

$$\mathbf{p}^{(n+1)} = -\mathbf{g}^{(n+1)} + \beta_n \mathbf{p}^{(n)}. \quad (10.41)$$

3. Метод сопряженных градиентов для минимизации неквадратичных функций. Для минимизации произвольной гладкой функции f расчетные формулы метода сопряженных градиентов модифицируют. Сначала полагают

$$\mathbf{g}^{(0)} = f'(\mathbf{x}^{(0)}), \quad \mathbf{p}^{(0)} = -\mathbf{g}^{(0)},$$

а затем при $n \geq 0$ вычисления производят по формулам

$$\varphi_n(\alpha_n) = \min_{\alpha \geq 0} \varphi_n(\alpha), \quad \text{где } \varphi_n(\alpha) = f(\mathbf{x}^{(n)} + \alpha \mathbf{p}^{(n)}). \quad (10.42)$$

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + \alpha_n \mathbf{p}^{(n)},$$

$$\begin{aligned}\mathbf{g}^{(n+1)} &= f'(\mathbf{x}^{(n+1)}), \\ \beta_n &= \frac{(\mathbf{g}^{(n+1)}, \mathbf{g}^{(n+1)})}{(\mathbf{g}^{(n)}, \mathbf{g}^{(n)})}, \\ \mathbf{p}^{(n+1)} &= -\mathbf{g}^{(n+1)} + \beta_n \mathbf{p}^{(n)}.\end{aligned}$$

Итерационный процесс здесь уже не оканчивается после конечного числа шагов, а направления $\mathbf{p}^{(n)}$ не являются, вообще говоря, сопряженными относительно некоторой матрицы.

Решение задач одномерной минимизации (10.42) приходится осуществлять численно. Отметим также то, что часто в методе сопряженных градиентов при $n = m, 2m, 3m, \dots$ коэффициент β_{n-1} не вычисляют, а полагают равным нулю. При этом очередной шаг производят фактически методом наискорейшего спуска. Такое «обновление» метода позволяет уменьшить влияние вычислительной погрешности.

Замечание. Иногда используется вариант метода, в котором коэффициент β_n вычисляют по формуле

$$\beta_n = \frac{(\mathbf{g}^{(n+1)}, \mathbf{g}^{(n+1)} - \mathbf{g}^{(n)})}{(\mathbf{g}^{(n)}, \mathbf{g}^{(n)})}.$$

Для сильно выпуклой гладкой функции f при некоторых дополнительных условиях метод сопряженных градиентов обладает высокой сверхлинейной скоростью сходимости. В то же время его трудоемкость невысока и сравнима с трудоемкостью метода наискорейшего спуска. Как показывает вычислительная практика, он незначительно уступает по эффективности квазиньютоновским методам, но предъявляет значительно меньшие требования к используемой памяти компьютера. В случае, когда решается задача минимизации функции с очень большим числом переменных, метод сопряженных градиентов, по-видимому, является единственным подходящим универсальным методом.

§ 10.6. Методы минимизации без вычисления производных

1. Методы прямого поиска. **Метод деформируемого многоугранника.** Существует большой класс методов минимизации, каждый из которых основан на сравнении значений целевой функции в последовательно вычисляемых пробных точках. Это так называемые *методы прямого поиска*. Обычно они применяются тогда, когда в любой окрестности точки локального минимума целевая функция не является гладкой,

а множество точек, в которых она недифференцируема, имеет слишком сложную структуру. К сожалению, методы прямого поиска в большинстве случаев очень неэффективны. Обращаться к ним, по-видимому, имеет смысл только тогда, когда есть уверенность, что никакие другие подходы к решению задачи минимизации невозможны.

Рассмотрим кратко один из наиболее известных методов прямого поиска — *метод деформируемого многогранника*. Сначала необходимо задать $m + 1$ точку $\mathbf{x}_1^{(0)}, \mathbf{x}_2^{(0)}, \dots, \mathbf{x}_{m+1}^{(0)}$ так, чтобы векторы $\mathbf{x}_2^{(0)} - \mathbf{x}_1^{(0)}$, $\mathbf{x}_3^{(0)} - \mathbf{x}_1^{(0)}$, ..., $\mathbf{x}_{m+1}^{(0)} - \mathbf{x}_1^{(0)}$ были линейно независимы. Тогда эти точки можно интерпретировать как вершины m -мерного многогранника. Обычно начальный многогранник выбирают правильным.

Опишем очередную $(n + 1)$ -ю итерацию простейшего варианта метода. Будем предполагать, что точки $\mathbf{x}_1^{(n)}, \mathbf{x}_2^{(n)}, \dots, \mathbf{x}_{m+1}^{(n)}$ являются вершинами многогранника, полученного на предыдущей итерации, и занумерованы так, что $f(\mathbf{x}_1^{(n)}) = \min_{1 \leq i \leq m+1} f(\mathbf{x}_i^{(n)})$, $f(\mathbf{x}_{m+1}^{(n)}) = \max_{1 \leq i \leq m+1} f(\mathbf{x}_i^{(n)})$.

Точку $\mathbf{x}_{m+1}^{(n)}$ отбрасывают как самую неудачную. Затем вычисляют «центр тяжести» оставшихся точек по формуле $\mathbf{u}^{(n)} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^{(n)}$ и строят пробную точку $\mathbf{x}_{\text{пп}}^{(n)} = \mathbf{u}^{(n)} + (\mathbf{u}^{(n)} - \mathbf{x}_{m+1}^{(n)})$, отразив симметрично вершину $\mathbf{x}_{m+1}^{(n)}$ относительно точки $\mathbf{u}^{(n)}$. Геометрическая иллюстрация этой операции для $m = 2$ приведена на рис. 10.9, а). Заметим, что в плоском случае роль многогранника играет треугольник.

Далее возможны три случая.

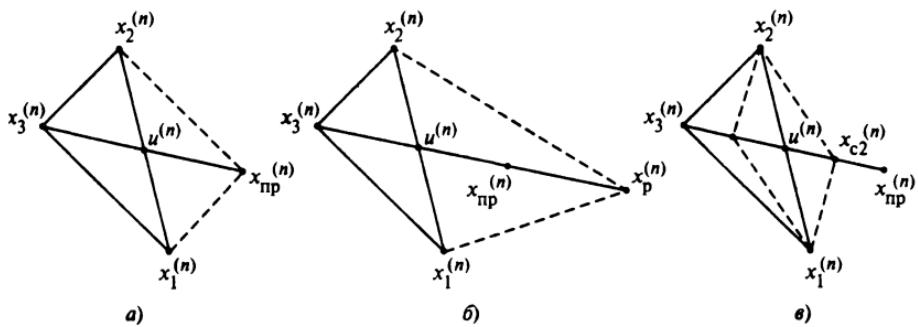


Рис. 10.9

Случай 1: $f(\mathbf{x}_1^{(n)}) \leq f(\mathbf{x}_{\text{пп}}^{(n)}) \leq f(\mathbf{x}_m^{(n)})$, т.е. пробная точка не привела к уменьшению значения целевой функции, но и не стала «худшей» в наборе $\mathbf{x}_1^{(n)}, \mathbf{x}_2^{(n)}, \dots, \mathbf{x}_m^{(n)}, \mathbf{x}_{\text{пп}}^{(n)}$. В этом случае точку признают удовлетворительной и $\mathbf{x}_{m+1}^{(n)}$ заменяют на $\mathbf{x}_{\text{пп}}^{(n)}$.

Случай 2: $f(\mathbf{x}_{\text{пп}}^{(n)}) < f(\mathbf{x}_1^{(n)})$, т.е. выбор пробной точки привел к уменьшению значения целевой функции. Тогда направление отражения считают удачным и делают попытку «растянуть» многогранник в этом направлении с коэффициентом растяжения $\beta > 1$ (обычно $2 \leq \beta \leq 3$). Находят точку $\mathbf{x}_p^{(n)} = \mathbf{u}^{(n)} + \beta(\mathbf{u}^{(n)} - \mathbf{x}_{m+1}^{(n)})$. Если $f(\mathbf{x}_p^{(n)}) < f(\mathbf{x}_{\text{пп}}^{(n)})$, то растяжение прошло удачно и точку $\mathbf{x}_{m+1}^{(n)}$, заменяют на $\mathbf{x}_p^{(n)}$. В противном случае растяжение было неудачным и производят замену $\mathbf{x}_{m+1}^{(n)}$ на $\mathbf{x}_{\text{пп}}^{(n)}$. На рис. 10.9, б показано «растяжение» с коэффициентом $\beta = 2$.

Случай 3: $f(\mathbf{x}_{\text{пп}}^{(n)}) > f(\mathbf{x}_m^{(n)})$. В этом случае считают, что многогранник следует сжать. Если $f(\mathbf{x}_{\text{пп}}^{(n)}) > f(\mathbf{x}_{m+1}^{(n)})$, то вершину $\mathbf{x}_{m+1}^{(n)}$ заменяют точкой $\mathbf{x}_{c_1}^{(n)} = \mathbf{u}^{(n)} - \gamma(\mathbf{u}^{(n)} - \mathbf{x}_{m+1}^{(n)})$. Если же $f(\mathbf{x}_{\text{пп}}^{(n)}) \leq f(\mathbf{x}_{m+1}^{(n)})$, то $\mathbf{x}_{m+1}^{(n)}$ заменяют на $\mathbf{x}_{c_2}^{(n)} = \mathbf{u}^{(n)} + \gamma(\mathbf{u}^{(n)} - \mathbf{x}_{m+1}^{(n)})$. Здесь $0 < \gamma < 1$, γ — коэффициент сжатия (обычно $\gamma = 1/2$). На рис. 10.9, в показано «сжатие» с коэффициентом $\gamma = 1/2$.

Помимо операций отражения, сжатия и растяжения периодически (после выполнения определенного числа итераций) производят операцию замены текущего многогранника правильным многогранником — так называемое *восстановление*. При восстановлении сохраняются лишь две лучшие точки последнего многогранника. Расстояние между этими точками принимают за длину ребра вновь генерируемого правильного многогранника. Эта операция позволяет ликвидировать излишние деформации, возникающие в ходе итераций.

Известно большое число различных модификаций этого метода. Подробно они описаны в книгах [28, 112]. Там же указаны различные критерии окончания.

2. Методы минимизации гладких функций, использующие конечно-разностные аппроксимации производных. Довольно часто в приложениях возникает необходимость в минимизации функций, обладающих достаточным числом производных, которые тем не менее

недоступны для прямого вычисления. Например, такая ситуация имеет место тогда, когда значение функции является результатом решения сложной математической задачи. В этом случае приходится применять алгоритмы, использующие лишь вычисляемые в различных точках значения функции f . Здесь обращение к методам прямого поиска (типа метода деформируемого многогранника), специально разработанным для минимизации негладких функций, вряд ли целесообразно, так как эти методы не позволяют извлекать никакой выгоды из возможной гладкости функции. Часто существенно лучший результат можно получить, заменив в одном из алгоритмов спуска используемые в нем производные их аппроксимациями в соответствии с формулами численного дифференцирования (см. гл. 12). Например, для производной f'_{x_j} простейшая аппроксимация такова:

$$f'_{x_j}(x_1, x_2, \dots, x_m) \approx \frac{1}{h_j} (f(x_1, \dots, x_{j-1}, x_j + h_j, x_{j+1}, \dots, x_m) - f(x_1, x_2, \dots, x_m)).$$

Подобная модификация алгоритмов не является тривиальной и требует достаточной осторожности в реализации. Это связано с высокой чувствительностью формул численного дифференцирования к ошибкам в вычислении функции (см. § 12.3). Более подробное обсуждение методов минимизации, использующих конечно-разностное аппроксимации производных, можно найти в [28] и [40]. В заключение все же отметим, что конечно-разностные аппроксимации производных используются в алгоритмах минимизации тогда, когда аналитическое вычисление производных невозможно. Если же производные можно вычислить аналитически, то усилия, потраченные на такое вычисление, как правило, окупаются.

§ 10.7. Нелинейная задача метода наименьших квадратов

1. Нелинейная задача метода наименьших квадратов. Среди задач безусловной минимизации особое место занимает задача минимизации функций вида

$$F(\mathbf{x}) = \frac{1}{2} \|\mathbf{f}(\mathbf{x})\|_2^2 = \frac{1}{2} \sum_{i=1}^N f_i(\mathbf{x})^2, \quad (10.43)$$

где $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_N(\mathbf{x}))^T$, $f_i(\mathbf{x})$ ($1 \leq i \leq N$) — некоторые функции m переменных $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$. Как правило, $N \geq m$; часто $N \gg m$.

Необходимость в минимизации таких функций возникает, например, при использовании метода наименьших квадратов для решения задачи аппроксимации функции (см. гл. 11) либо для решения проблемы идентификации математической модели по данным эксперимента.

К поиску точки минимума функции вида (10.43) может быть сведена и задача решения системы нелинейных уравнений

$$f_i(x) = 0, \quad 1 \leq i \leq N,$$

то есть системы

$$f(x) = 0. \quad (10.44)$$

Действительно, если решение \bar{x} системы (10.44) существует, то оно совпадает с точкой глобального минимума функции (10.43), причем $F(\bar{x}) = 0$. В случае, когда число уравнений в системе превышает число неизвестных, мы имеем дело с переопределенной системой, которая, скорее всего не имеет решений. В этом случае минимизация функции F позволяет найти вектор \bar{x} , для которого

$$f(\bar{x}) \approx 0$$

с минимальной нормой невязки $f(\bar{x})$.

Для минимизации функции (10.43), в принципе, можно применить любой из универсальных методов минимизации. Однако, как правило, этого не делают, а используют алгоритмы, разработанные специально для решения нелинейной задачи метода наименьших квадратов. В частности, это связано с тем, что градиент $g(x) = F'(x)$ и матрица Гессе $G(x) = F''(x)$ для функции (10.43) имеют специальный вид:

$$g(x) = J(x)^T f(x); \quad (10.45)$$

$$G(x) = J(x)^T J(x) + Q(x), \text{ где } Q(x) = \sum_{i=1}^N f_i(x) G_i(x). \quad (10.46)$$

Здесь $J(x) = f'(x)$ — это матрица Якоби для вектор-функции f с элементами $J_{ij} = \frac{\partial f_i}{\partial x_j}(x)$, $1 \leq i \leq N$, $1 \leq j \leq m$.¹ Матрицы $G_i(x)$ ($1 \leq i \leq N$) —

это матрицы Гессе для функций f_i с элементами $(G_i)_{kj} = \frac{\partial^2 f_i}{\partial x_k \partial x_j}(x)$, $1 \leq k \leq m$, $1 \leq j \leq m$.

2. Метод Гаусса—Ньютона. Будем рассматривать задачу минимизации функции (10.43) как задачу решения переопределенной системы

¹ Обратим внимание на то, что в отличие от матрицы Гессе (7.3) матрица $J(x)$ уже не квадратная, а прямоугольная.

(10.44). Пусть некоторое приближение $\mathbf{x}^{(k)}$ к точке минимума найдено. Аналогично тому, как мы поступали в § 7.3, произведем линеаризацию вектор-функции \mathbf{f} в этой точке, заменив ее главной линейной частью разложения по формуле Тейлора в точке $\mathbf{x}^{(k)}$:

$$\mathbf{f}(\mathbf{x}) \approx \mathbf{f}(\mathbf{x}^{(k)}) + \mathbf{J}^{(k)}(\mathbf{x} - \mathbf{x}^{(k)});$$

здесь и всюду в этом параграфе $\mathbf{J}^{(k)} = \mathbf{J}(\mathbf{x}^{(k)}) = \mathbf{f}'(\mathbf{x}^{(k)})$. За следующее приближение $\mathbf{x}^{(k+1)}$ к решению примем точку минимума функции

$$F_k(\mathbf{x}) = \frac{1}{2} \|\mathbf{f}(\mathbf{x}^{(k)}) + \mathbf{J}^{(k)}(\mathbf{x} - \mathbf{x}^{(k)})\|_2^2.$$

Для этого положим $\mathbf{p} = \mathbf{x} - \mathbf{x}^{(k)}$ и найдем вектор $\mathbf{p}^{(k+1)}$ такой, что

$$\|\mathbf{f}(\mathbf{x}^{(k)}) + \mathbf{J}^{(k)}\mathbf{p}^{(k+1)}\|_2^2 = \min_{\mathbf{p} \in \mathbb{R}^m} \|\mathbf{f}(\mathbf{x}^{(k)}) + \mathbf{J}^{(k)}\mathbf{p}\|_2^2. \quad (10.47)$$

Если решение этой задачи не является единственным, то за $\mathbf{p}^{(k+1)}$ принимается решение с минимальной нормой.

После того, как $\mathbf{p}^{(k+1)}$ найдено, приближение $\mathbf{x}^{(k+1)}$ вычисляется по формуле

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^{(k+1)}.$$

Заметим, что задача (10.47) — это линейная задача метода наименьших квадратов, предназначенная для решения переопределенной системы

$$\mathbf{J}(\mathbf{x}^{(k)})\mathbf{p} = -\mathbf{f}(\mathbf{x}^{(n)}).$$

Ее решение $\mathbf{p}^{(k+1)}$ может быть найдено одним из рассмотренных в § 5.12 методов.

Замечание 1. Можно показать, что в случае, когда решение $\mathbf{p}^{(k+1)}$ задачи (10.47) ненулевое, оно является направлением спуска для функции F .

Замечание 2. В случае, когда $m = N$, а матрица Якоби $\mathbf{J}(\mathbf{x})$ невырождена, метод Гаусса—Ньютона превращается в метод Ньютона, рассмотренный в § 7.3.

3. Модифицированный метод Гаусса—Ньютона. Как и метод Ньютона, метод Гаусса—Ньютона обладает лишь локальной сходимостью. Для того, чтобы расширить область сходимости метода его модифицируют.

После того, как решение линейной задачи метода наименьших квадратов (10.47) найдено, полагают

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k+1)},$$

где $\alpha_k > 0$ вычисляется как решение задачи одномерной минимизации

$$\varphi_k(\alpha_k) = \min_{\alpha \geq 0} \varphi_k(\alpha), \quad \text{где } \varphi_k(\alpha) = \|f(x^{(k)}) + \alpha p^{(k+1)}\|_2^2.$$

4. Другие методы. Коротко обсудим другие подходы к решению нелинейной задачи метода наименьших квадратов.

Заметим, что искомое решение должно удовлетворять необходимому условию экстремума $F'(x) = 0$. В силу формулы (10.45) это условие принимает вид

$$J(x)^T f(x) = 0. \quad (10.47)$$

Применяя для решения этой системы нелинейных уравнений классический метод Ньютона, мы приходим к следующему итерационному методу

$$G^{(k)} p^{(k+1)} = -[J^{(k)}]^T f(x^{(k)}), \quad (10.48)$$

$$x^{(k+1)} = x^{(k)} + p^{(k+1)}. \quad (10.49)$$

Матрица $G^{(k)} = G(x^{(k)})$, которую следовало бы использовать в этом методе, весьма громоздка и неудобна для вычислений. Согласно формуле (10.46) она имеет вид

$$G^{(k)} = [J^{(k)}]^T J^{(k)} + Q^{(k)}, \quad \text{где } Q^{(k)} = \sum_{i=1}^N f_i(x^{(k)}) G_i(x^{(k)}).$$

Для вычисления слагаемого $Q^{(k)}$ на каждой итерации требуется вычисление всех вторых частных производных всех функций f_i , $1 \leq i \leq N$. Этого всячески стараются избегать.

Обратим внимание на то, что если получено достаточно хорошее приближение $x^{(k)}$ такое, что $\|f(x^{(k)})\|_2 \approx 0$, то $f_i(x^{(k)}) \approx 0$ для всех $i = 1, 2, \dots, N$. В этом случае $Q^{(k)} \approx 0$ и $G^{(k)} \approx [J^{(k)}]^T J^{(k)}$.

Замена в итерационном методе (10.48), (10.49) матрицы $G^{(k)}$ на матрицу $[J^{(k)}]^T J^{(k)}$ приводит к итерационному методу

$$[J^{(k)}]^T J^{(k)} p^{(k+1)} = -J(x^{(k)})^T f(x^{(k)}); \quad (10.50)$$

$$x^{(k+1)} = x^{(k)} + p^{(k+1)}, \quad (10.51)$$

который эквивалентен методу Гаусса—Ньютона, так как решения системы (10.50) совпадают с точками минимума функции

$$\|f(x^{(k)}) + J(x^{(k)})p\|_2^2.$$

Распространенной альтернативой методу Гаусса—Зейделя является метод Левенберга—Маркардта. В этом методе матрица $G^{(k)}$ заменяется

матрицей $[\mathbf{J}^{(k)}]^T \mathbf{J}^{(k)} + \alpha_k \mathbf{E}$, где $0 \leq \alpha_k$ — параметр, и метод принимает вид

$$([\mathbf{J}^{(k)}]^T \mathbf{J}^{(k)} + \alpha_k \mathbf{E}) \mathbf{p}^{(k+1)} = -[\mathbf{J}^{(k)}]^T \mathbf{f}(\mathbf{x}^{(k)}),$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^{(k+1)}.$$

Методы Гаусса—Ньютона и Левенберга—Маркардта исходят из того, что вблизи от решения невязка $\mathbf{f}(\mathbf{x}^{(k)})$ мала и поэтому слагаемым $\mathbf{Q}^{(k)}$ в формуле для вычисления $\mathbf{G}^{(k)}$ можно пренебречь. Такие задачи обычно называют задачами с малыми невязками.

Однако на практике часто встречаются задачи с большими невязками, в которых это слагаемое существенно. В этом случае строят специальные квазиньютоновские приближения $\mathbf{M}^{(k)}$ к матрице $\mathbf{Q}^{(k)}$ и итерации проводят по формулам

$$([\mathbf{J}^{(k)}]^T \mathbf{J}^{(k)} + \mathbf{M}^{(k)}) \mathbf{p}^{(k+1)} = -\mathbf{J}(\mathbf{x}^{(k)})^T \mathbf{f}(\mathbf{x}^{(k)}),$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^{(k+1)}.$$

Достаточно подробное обсуждение квазиньютоновских методов и других методов решения нелинейной задачи метода наименьших квадратов можно найти в книгах [28, 40].

§ 10.8. Дополнительные замечания

1. Естественно, что в данной книге содержится лишь краткое введение в методы решения задач безусловной минимизации. Более подробное изложение можно найти, например в [6, 20, 28, 40, 78, 80, 94, 112 и др.].

2. Задачи условной минимизации. Множество X , на котором минимизируется функция f , часто задается с помощью системы неравенств вида

$$g_i(x_1, x_2, \dots, x_m) \geq 0, \quad i = 1, 2, \dots, k.$$

Если целевая функция f и задающие множество X функции g_i линейны, то задачу минимизации называют *задачей линейного программирования*. Если же хотя бы одна из этих функций нелинейна, то говорят о *задаче нелинейного программирования*.

Обсуждение методов решения задач условной минимизации при всей их важности выходит за рамки данной книги. Укажем, например на книги [20, 28, 94, 112], содержащие достаточно подробное и доступное введение в соответствующие методы.

3. Особый класс задач составляют так называемые *задачи дискретной минимизации*. В этих задачах множество X , на котором миними-

зируется функция f , является конечным или счетным. Часто X — множество точек с целочисленными координатами, удовлетворяющими некоторым ограничениям. Методы решения таких задач описаны, например в [94].

4. Авторы рекомендуют обратить внимание на книгу [28], в особенности на ее последние главы («Моделирование», «Практические вопросы» и «Вопросы и ответы»). В них содержится большое число весьма интересных замечаний, полезных рекомендаций и советов для широкого круга специалистов, заинтересованных в решении практических задач минимизации.

ПРИБЛИЖЕНИЕ ФУНКЦИЙ И СМЕЖНЫЕ ВОПРОСЫ

В этой главе рассматриваются наиболее важные и часто встречающиеся в приложениях методы приближения функций одной переменной. Значительное внимание уделено интерполяции, причем рассматривается интерполяция не только алгебраическими многочленами, но и тригонометрическими многочленами, а также интерполяция сплайнами. Довольно подробно обсуждается метод наименьших квадратов, широко используемый в практике инженерных расчетов. Дается понятие о наилучшем равномерном приближении и дробно-рациональных аппроксимациях.

В главу включены также некоторые вопросы вычислительной математики, имеющие непосредственное отношение к методам аппроксимации (приближения) функций. Это конечные и разделенные разности, многочлены Чебышева, быстрое дискретное преобразование Фурье.

§ 11.1. Постановка задачи приближения функций

Вычисление значения функции $y = f(x)$ — одна из тех задач, с которой на практике постоянно приходится сталкиваться. Естественно, что при решении на компьютере серьезных задач желательно иметь быстрые и надежные алгоритмы вычисления значений используемых функций. Для элементарных, а также для основных специальных функций такие алгоритмы разработаны, реализованы в виде стандартных программ и включены в математическое обеспечение компьютера. Однако в расчетах нередко используются и другие функции, непосредственное вычисление которых затруднено либо приводит к слишком большим затратам машинного времени. Укажем на некоторые типичные ситуации.

1. Функция f задана таблицей своих значений:

$$y_i = f(x_i), \quad i = 0, 1, 2, \dots, n, \quad (11.1)$$

а вычисления производятся в точках x , не совпадающих с табличными.

2. Непосредственное вычисление значения $y = f(x)$ связано с проведением сложных расчетов и приводит к значительным затратам машинного времени, которые могут оказаться неприемлемыми, если функция f вычисляется многократно.

3. При заданном значении x значение $f(x)$ может быть найдено из эксперимента. Ясно, что такой способ «вычисления» в большинстве случаев нельзя использовать в вычислительных алгоритмах, так как он связан с необходимостью прерывания вычислительного процесса для проведения эксперимента¹. В этой ситуации экспериментальные данные получают до начала вычислений на компьютере. Нередко они представляют собой таблицу типа (11.1) с тем отличием, что табличные значения y_i^* отличаются от «истинных» значений y_i , так как заведомо содержат ошибки эксперимента.

Возникающие проблемы нередко удается решить следующим образом. Функцию $f(x)$ приближенно заменяют другой функцией $g(x)$, вычисляемые значения которой и принимают за приближенные значения функции f . Конечно, такая замена оправдана лишь тогда, когда значения $g(x)$ вычисляются быстро и надежно, а погрешность приближения $f(x) - g(x)$ достаточно мала. Обсудим кратко некоторые вопросы, с которыми в каждом конкретном случае приходится сталкиваться при выборе постановки задачи приближения и метода ее решения.

1°. Необходимо решить, какую информацию о функции f можно использовать как входные данные для вычисления приближения g . Например, часто известна или может быть получена таблица значений функции вида (11.1), а иногда — и таблица ее производных. В некоторых случаях можно использовать информацию о значениях функции на всем отрезке $[a, b]$.

2°. Полезно иметь некоторую дополнительную априорную информацию об аппроксимируемой функции. Часто она бывает качественного характера, например известно, что функция f «достаточно гладкая» («плавно меняющаяся»), периодическая, монотонная, четная и т.п. Иногда удается получить некоторые количественные характеристики функции f , например, бывают известны верхние оценки для максимума модуля некоторых ее производных, величина периода, оценка уровня погрешности в заданных значениях.

3°. Знание свойств функции f позволяет осознанно выбирать класс G аппроксимирующих функций. Часто такой класс представляет собой параметрическое семейство функций вида $y = g(x, \alpha) = g(x, a_0, a_1, \dots, a_m)$ и выбор конкретной аппроксимирующей функции g осуществляется с помощью выбора параметров a_0, a_1, \dots, a_m .

¹Правда, в некоторых алгоритмах такое прерывание естественно, например, если компьютер используется для управления технологическим процессом, сложной технической системой или включен в систему обработки и планирования физического эксперимента.

Широко используются классы функций вида

$$\Phi_m(x) = a_0 \varphi_0(x) + a_1 \varphi_1(x) + \dots + a_m \varphi_m(x), \quad (11.2)$$

являющихся линейными комбинациями фиксированного набора некоторых базисных функций $\varphi_0(x)$, $\varphi_1(x)$, ..., $\varphi_m(x)$. Функцию $\Phi_m(x)$ часто называют *обобщенным многочленом* по системе функций φ_0 , φ_1 , ..., φ_m , а число m — его *степенью*.

Если в качестве базисных функций берутся степенные функции $\varphi_k(x) = x^k$, то возникает задача приближения алгебраическими многочленами

$$P_m(x) = a_0 + a_1 x + \dots + a_m x^m. \quad (11.3)$$

Отметим, что методы приближения функций алгебраическими многочленами играют важную роль в численном анализе и наиболее глубоко разработаны. Одна из причин этого состоит в том, что многочлены (11.3) легко вычисляются, без труда дифференцируются и интегрируются.

Тригонометрические многочлены

$$S_m(x) = \alpha_0 + \sum_{1 \leq k \leq m/2} (\alpha_k \cos 2\pi kx + \beta_k \sin 2\pi kx), \quad (11.4)$$

часто используемые для аппроксимации периодических на отрезке $[0, 1]$ функций, также могут быть записаны в виде (11.2), если в качестве базисных функций выбрать функции $\varphi_0(x) = 1$, $\varphi_1(x) = \cos 2\pi x$, $\varphi_2(x) = \sin 2\pi x$, $\varphi_3(x) = \cos 4\pi x$, $\varphi_4(x) = \sin 4\pi x$, Используя формулу Эйлера $\exp\{iy\} = \cos y + i \sin y$, запишем тригонометрический многочлен (11.4) в виде

$$S_m(x) = \sum_{-m/2 \leq k \leq m/2} a_k \exp(2\pi ikx), \quad (11.5)$$

что соответствует выбору базисных функций $\varphi_k(x) = \exp(2\pi ikx)$, $-m/2 \leq k \leq m/2$.

Используются также и некоторые нелинейные комбинации функций, отличные от (11.2). Например, в ряде случаев эффективным является использование класса дробно-рациональных функций

$$\frac{a_0 + a_1 x + \dots + a_m x^m}{1 + b_1 x + \dots + b_k x^k}.$$

Выбор класса G аппроксимирующих функций осуществляется с учетом того, насколько хорошо может быть приближена функция f функциями из этого класса.

4°. Необходим критерий выбора в классе G конкретной аппроксимирующей функции g , являющейся в смысле этого критерия наилучшим

приближением к f . Например, требование совпадения функции g с функцией f в некоторых фиксированных точках приводит к задаче интерполяции. Другой распространенный критерий — требование минимизации среднеквадратичного уклонения — лежит в основе метода наименьших квадратов. Существует большое число других критериев, естественных в конкретных прикладных проблемах.

5°. Важно понимать, что решение указанных выше вопросов тесно связано с тем, как мы собираемся использовать приближение g и какая точность нам нужна.

Замечание. Задачу выбора в классе G конкретной приближающей функции можно рассматривать как задачу идентификации (см. § 1.1), если интерпретировать функцию $y = g(x, a)$ как математическую модель реальной функциональной зависимости $y = f(x)$.

§ 11.2. Интерполяция обобщенными многочленами

1. Постановка задачи интерполяции. Пусть в точках x_0, x_1, \dots, x_n , расположенных на отрезке $[a, b]$ и попарно различных, задана таблица (11.1) значений некоторой функции f . Задача интерполяции состоит в построении функции g , удовлетворяющей условию

$$g(x_i) = y_i \quad (i = 0, 1, \dots, n). \quad (11.6)$$

Другими словами, ставится задача о построении функции g , график которой проходит через заданные точки (x_i, y_i) (рис. 11.1). Указанный способ приближения функций принято называть *интерполяцией* (или *интерполированием*), а точки x_i — *узлами интерполяции*.

Нетрудно видеть, что выбор функции g неоднозначен, так как по заданной таблице можно построить бесконечно много интерполирующих функций. На практике, как правило, функцию g выбирают из достаточно узкого класса G функций, в котором единственность выбора гарантируется.

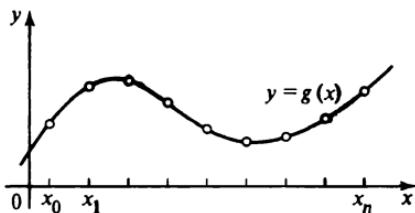


Рис. 11.1

2. Экстраполяция. Пусть x_{\min} и x_{\max} — минимальный и максимальный из узлов интерполяции. В случае, когда интерполяция используется для вычисления приближенного значения функции f в точке x , не принадлежащей отрезку $[x_{\min}, x_{\max}]$ (*отрезку наблюдения*), принято говорить о том, что осуществляется **экстраполяция**. Этот метод приближения часто используют в целях прогнозирования характера протекания тех или иных процессов при значениях параметров x , выходящих за пределы отрезка наблюдения. Заметим, что надежность такого прогноза при значениях x , удаленных на значительное расстояние от отрезка $[x_{\min}, x_{\max}]$, как правило, невелика.

3. Задача интерполяции обобщенными многочленами. Рассмотрим более подробно задачу интерполяции обобщенными многочленами $\Phi_m(x)$ вида (11.2). Назовем обобщенный многочлен $\Phi_m(x)$ **интерполяционным**, если он удовлетворяет условию

$$\Phi_m(x_i) = y_i \quad (i = 0, 1, \dots, n), \quad (11.6)$$

или, что то же самое, системе линейных алгебраических уравнений

$$\begin{aligned} \varphi_0(x_0)a_0 + \varphi_1(x_0)a_1 + \dots + \varphi_m(x_0)a_m &= y_0, \\ \varphi_0(x_1)a_0 + \varphi_1(x_1)a_1 + \dots + \varphi_m(x_1)a_m &= y_1, \\ \dots \dots \dots & \\ \varphi_0(x_n)a_0 + \varphi_1(x_n)a_1 + \dots + \varphi_m(x_n)a_m &= y_n \end{aligned} \quad (11.7)$$

относительно коэффициентов a_0, a_1, \dots, a_m .

Заметим, что систему уравнений (11.7) можно записать в следующем виде:

$$\mathbf{P}\mathbf{a} = \mathbf{y}, \quad (11.8)$$

где

$$\mathbf{P} = \begin{bmatrix} \varphi_0(x_0) & \varphi_1(x_0) & \dots & \varphi_m(x_0) \\ \varphi_0(x_1) & \varphi_1(x_1) & \dots & \varphi_m(x_1) \\ \dots \dots \dots \\ \varphi_0(x_n) & \varphi_1(x_n) & \dots & \varphi_m(x_n) \end{bmatrix}, \quad \mathbf{a} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_m \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}. \quad (11.9)$$

Введем векторы $\varphi_j = (\varphi_j(x_0), \varphi_j(x_1), \dots, \varphi_j(x_n))^T, j = 0, 1, \dots, m$. Будем говорить, что *система функций $\varphi_0, \varphi_1, \dots, \varphi_m$ линейно зависима в точках x_0, x_1, \dots, x_n* , если один из векторов φ_j системы $\varphi_0, \varphi_1, \dots, \varphi_m$ может быть

представлен в виде линейной комбинации остальных векторов этой системы:

$$\varphi_j = \sum_{\substack{k=0 \\ k \neq j}}^m \alpha_k \varphi_k. \quad (11.10)$$

В противном случае систему функций $\varphi_0, \varphi_1, \dots, \varphi_m$ будем называть *линейно независимой в точках x_0, x_1, \dots, x_n* .

Пример 11.1. Покажем, что при $m \leq n$ система функций $1, x, x^2, \dots, x^m$ линейно независима в точках x_0, x_1, \dots, x_n , если они попарно различны.

Допустим противное. Тогда справедливо равенство (11.10), которое в данном случае при $\varphi_k = (x_0^k, x_1^k, \dots, x_n^k)^T$ принимает вид

$$x_i^j = \sum_{\substack{k=0 \\ k \neq j}}^m \alpha_k x_i^k \quad (i=0, 1, \dots, n). \quad (11.11)$$

Полагая $\alpha_j = -1$, получаем, что многочлен $P_m(x) = \sum_{k=0}^m \alpha_k x^k$ степени¹ m

обращается в ноль в точках x_0, x_1, \dots, x_n , число которых равно $n+1$ и, следовательно, больше m . Однако в силу основной теоремы алгебры многочлен степени m , тождественно не равный нулю, не может иметь более m корней. Полученное противоречие доказывает линейную независимость рассматриваемой системы функций. ▲

Рассмотрим *матрицу Грама* системы функций $\varphi_0, \varphi_1, \dots, \varphi_m$, имеющую вид

$$\Gamma = P^* P = \begin{bmatrix} (\varphi_0, \varphi_0) & (\varphi_1, \varphi_0) & \dots & (\varphi_m, \varphi_0) \\ (\varphi_0, \varphi_1) & (\varphi_1, \varphi_1) & \dots & (\varphi_m, \varphi_1) \\ \dots & \dots & \dots & \dots \\ (\varphi_0, \varphi_m) & (\varphi_1, \varphi_m) & \dots & (\varphi_m, \varphi_m) \end{bmatrix}. \quad (11.12)$$

¹ В данной главе для упрощения формулировок будем говорить, что многочлен P_m имеет степень m даже в случае, когда $a_m = 0$, т.е. фактическая его степень меньше m .

Здесь в случае, когда функции φ_j могут принимать комплексные значения, под P^* понимается сопряженная к P матрица, а элементы γ_{jk} матрицы Грама вычисляются по формуле

$$\gamma_{jk} = (\varphi_k, \varphi_j) = \sum_{i=0}^n \varphi_k(x_i) \overline{\varphi_j(x_i)}. \quad (11.13)$$

Если же функции φ_j принимают только вещественные значения, то $P^* = P^T$ и элементы матрицы Грама вычисляются по формуле

$$\gamma_{jk} = (\varphi_k, \varphi_j) = \sum_{i=0}^n \varphi_k(x_i) \varphi_j(x_i). \quad (11.14)$$

Определитель матрицы Грама $\det \Gamma$ принято называть *определителем Грама*.

Как следует из курса линейной алгебры, справедливо следующий результат.

Теорема 11.1. Система функций $\varphi_0, \varphi_1, \dots, \varphi_m$ является линейно независимой в точках x_0, x_1, \dots, x_m тогда и только тогда, когда $m \leq n$ и определитель Грама $\det \Gamma$ отличен от нуля. ■

Известно, что при $m > n$ система функций $\varphi_0, \varphi_1, \dots, \varphi_m$ линейно зависима в точках x_0, x_1, \dots, x_n . Отсюда вытекает неединственность решения α системы (11.8) (если оно существует). Действительно, в этом случае справедливо представление (11.10) и вместе с вектором α решением системы (11.8) является вектор $\alpha' = \alpha + t\Delta\alpha$, где $\Delta\alpha = (\alpha_0, \alpha_1, \dots, \alpha_{j-1}, -1, \alpha_{j+1}, \dots, \alpha_m)^T$, а t — любое число. Если же $m < n$, то решение системы (11.8) существует не для всякой правой части y .

В силу указанных причин при интерполяции обобщенными многочленами число параметров $m + 1$ обычно берут равным числу $n + 1$ заданных точек. В этом случае P — квадратная матрица и для того, чтобы система (11.8) была однозначно разрешима при любой правой части y , необходимо и достаточно, чтобы определитель матрицы P был отличен от нуля. В свою очередь при $m = n$ это условие в силу равенства $\det \Gamma = \det P^* \det P = |\det P|^2$ и теоремы 11.1 дает следующий результат.

Теорема 11.2. Если $m = n$, то решение задачи интерполяции обобщенным многочленом (11.2) существует и единственно при любом наборе данных y_0, y_1, \dots, y_n тогда и только тогда когда система функций $\varphi_0, \varphi_1, \dots, \varphi_n$ линейно независима в точках x_0, x_1, \dots, x_n . ■

Назовем систему функций $\varphi_0, \varphi_1, \dots, \varphi_m$ ортогональной на множестве точек x_0, x_1, \dots, x_n , если $(\varphi_k, \varphi_j) = 0$ при $k \neq j$ и $(\varphi_k, \varphi_j) \neq 0$ при $k = j$ для всех $k = 0, 1, \dots, m; j = 0, 1, \dots, m$. Очевидно, что для ортогональной на множестве x_0, x_1, \dots, x_n системы функций матрица Грама диагональна, а определитель Грама отличен от нуля. Поэтому всякая ортогональная на множестве точек x_0, x_1, \dots, x_n система функций заведомо является линейно независимой в этих точках.

Пример 11.2. Покажем, что система функций $\varphi_0, \varphi_1, \dots, \varphi_{N-1}$, где $\varphi_k(x) = \exp\{2\pi ikx\}$, ортогональна на множестве точек $x_l = l/N, l = 0, 1, \dots, N-1$. Здесь i — мнимая единица.

Для доказательства ортогональности рассматриваемой системы функций достаточно установить, что справедливо равенство

$$(\varphi_k, \varphi_j) = N\delta_{kj}, \quad 0 \leq k < N, 0 \leq j < N, \quad (11.15)$$

где $\delta_{kj} = 0$ при $k \neq j$ и $\delta_{kj} = 1$ при $k = j$. Введем обозначение $\omega = \exp\{2\pi i/N\}$. Тогда $\varphi_k(x_l) = \exp\{2\pi ikl/N\} = \omega^{kl}$ и согласно формуле (11.13) имеем

$$(\varphi_k, \varphi_j) = \sum_{l=0}^{N-1} \omega^{kl} \omega^{-jl} = \sum_{l=0}^{N-1} [\omega^{k-j}]^l. \quad (11.16)$$

При $k = j$ правая часть равенства (11.16), очевидно, равна N . При $k \neq j$, используя формулу суммы членов геометрической прогрессии и равенство $\omega^{(k-j)N} = \exp\{2\pi i(k-j)\} = 1$, имеем

$$(\varphi_k, \varphi_j) = (1 - [\omega^{k-j}]^N)/(1 - \omega^{k-j}) = 0.$$

Таким образом, равенство (11.15), а вместе с ним и ортогональность системы функций $\varphi_k(x) = \exp\{2\pi ikx\}$ доказаны. ▲

Когда система функций $\varphi_0, \varphi_1, \dots, \varphi_n$ ортогональна на множестве точек x_0, x_1, \dots, x_n , решение задачи интерполяции не представляет затруднений. Действительно, система уравнений (11.8) после умножения на матрицу P^* преобразуется к виду

$$\Gamma \mathbf{a} = \mathbf{b}, \text{ где } \mathbf{b} = P^* \mathbf{y}. \quad (11.17)$$

Заметим, что элементы вектора $\mathbf{b} = (b_0, b_1, \dots, b_m)^T$ вычисляются по формуле

$$b_j = (\mathbf{y}, \varphi_j) = \sum_{l=0}^n y_l \overline{\varphi_j(x_l)}, \quad j = 0, 1, \dots, m. \quad (11.18)$$

Так как матрица Γ диагональна, то решение системы (11.17) находится в явном виде:

$$a_j = \frac{(y, \varphi_j)}{(\varphi_j, \varphi_j)}, \quad j = 0, 1, \dots, m. \quad (11.19)$$

§ 11.3. Полиномиальная интерполяция. Многочлен Лагранжа

1. Интерполяционный многочлен. Начнем с рассмотрения задачи интерполяции в наиболее простом и полно исследованном случае интерполирования алгебраическими многочленами. Для заданной таблицы (11.1) многочлен $P_n(x) = \sum_{k=0}^n a_k x^k$ степени n называется *интерполяционным многочленом*, если он удовлетворяет условиям

$$P_n(x_i) = y_i, \quad i = 0, 1, \dots, n. \quad (11.20)$$

Равенства (11.20) можно записать аналогично (11.7) в виде системы уравнений

$$\begin{aligned} a_0 + a_1 x_0 + a_2 x_0^2 + \dots + a_n x_0^n &= y_0, \\ a_0 + a_1 x_1 + a_2 x_1^2 + \dots + a_n x_1^n &= y_1, \\ \dots & \\ a_0 + a_1 x_n + a_2 x_n^2 + \dots + a_n x_n^n &= y_n \end{aligned} \quad (11.21)$$

относительно коэффициентов многочлена. Эта система однозначно разрешима, так как система функций $1, x, x^2, \dots, x^n$ линейно независима в точках x_0, x_1, \dots, x_n (см. пример 11.1 и теорему 11.2). Однозначная разрешимость системы (11.21) следует и из того хорошо известного факта, что определитель этой системы (определитель Вандермонда¹)

$$\begin{vmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{vmatrix} = \prod_{0 \leq j < i \leq n} (x_i - x_j)$$

¹ Александр Теофил Вандермонд (1735—1796) — французский математик.

отличен от нуля, если узлы интерполяции попарно различны. Таким образом, верна следующая теорема.

Теорема 11.3. Существует единственный интерполяционный многочлен степени n , удовлетворяющий условиям (11.20). ■

Замечание. На практике система (11.21) никогда не используется для вычисления коэффициентов интерполяционного многочлена. Дело в том, что часто она является плохо обусловленной. Кроме того, существуют различные удобные явные формы записи интерполяционного многочлена, которые и применяются при интерполяции. Наконец, в большинстве приложений интерполяционного многочлена явное вычисление коэффициентов a_k не нужно.

2. Многочлен Лагранжа. Приведем одну из форм записи интерполяционного многочлена — **многочлен Лагранжа**

$$L_n(x) = \sum_{j=0}^n y_j l_{nj}(x). \quad (11.22)$$

Здесь

$$l_{nj}(x) = \prod_{\substack{k=0 \\ k \neq j}}^n \frac{x - x_k}{x_j - x_k} = \frac{(x - x_0)(x - x_1) \dots (x - x_{j-1})(x - x_{j+1}) \dots (x - x_n)}{(x_j - x_0)(x_j - x_1) \dots (x_j - x_{j-1})(x_j - x_{j+1}) \dots (x_j - x_n)}.$$

Как нетрудно видеть, $l_{nj}(x)$ представляет собой многочлен степени n , удовлетворяющий условию

$$l_{nj}(x_i) = \begin{cases} 1 & \text{при } i=j \\ 0 & \text{при } i \neq j \end{cases}.$$

Таким образом, степень многочлена L_n равна n и при $x = x_i$ в сумме (11.22) обращаются в ноль все слагаемые, кроме слагаемого с номером $j = i$, равного y_i . Поэтому многочлен Лагранжа (11.22) действительно является интерполяционным.

Замечание 1. Запись интерполяционного многочлена в форме Лагранжа (11.22) можно рассматривать как его запись в виде обобщенного многочлена (11.2) по системе функций $\varphi_k(x) = l_{nk}(x)$, $k = 0, 1, \dots, n$.

Замечание 2. Как правило, интерполяционный многочлен Лагранжа используется так, что нет необходимости его преобразования к каноническому виду $L_n(x) = \sum_{k=0}^n a_k x^k$. Более того, часто такое преобразование нежелательно.

В инженерной практике наиболее часто используется интерполяция многочленами первой, второй и третьей степени (*линейная, квадратичная и кубическая интерполяции*). Приведем соответствующие формулы для записи многочленов Лагранжа первой и второй степени:

$$L_1(x) = y_0 \frac{x - x_1}{x_0 - x_1} + y_1 \frac{x - x_0}{x_1 - x_0}, \quad (11.23)$$

$$L_2(x) = y_0 \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} + y_1 \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} + y_2 \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}. \quad (11.24)$$

Пример 11.3. Пусть задана таблица значений функции $y = \ln x$ (табл. 11.1).

Таблица 11.1

x	1.0	1.1	1.2	1.3	1.4
y	0.000000	0.095310	0.182322	0.262364	0.336472

Для приближенного вычисления значения $\ln(1.23)$ воспользуемся линейной и квадратичной интерполяцией.

Возьмем $x_0 = 1.2$, $x_1 = 1.3$. Вычисление по формуле (11.23) дает значение $\ln(1.23) \approx L_1(1.23) \approx 0.206335$.

Для применения квадратичной интерполяции возьмем $x_0 = 1.1$, $x_1 = 1.2$, $x_2 = 1.3$ — три ближайших к точке $x = 1.23$ узла. Вычисляя по формуле (11.24), имеем $\ln(1.23) \approx L_2(1.23) \approx 0.207066$.

Заметим, что пока нам не известна погрешность полученных приближенных значений. ▲

§ 11.4. Погрешность интерполяции

Приведем без доказательства наиболее известную теорему о погрешности интерполяции.

Теорема 11.4. Пусть функция f дифференцируема $n + 1$ раз на отрезке $[a, b]$, содержащем узлы интерполяции x_i , $i = 0, 1, \dots, n$. Тогда для погрешности интерполяции в точке $x \in [a, b]$ справедливо равенство

$$f(x) - P_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{n+1}(x), \quad (11.25)$$

в котором $\omega_{n+1}(x) = (x - x_0)(x - x_1) \dots (x - x_n)$, а ξ — некоторая точка, принадлежащая интервалу (a, b) . ■

Основное неудобство в использовании этой теоремы состоит в том, что входящая в формулу (11.25) для погрешности точка ξ неизвестна. Поэтому чаще используется не сама теорема, а ее следствие.

Следствие. В условиях теоремы 11.4 справедлива оценка погрешности интерполяции в точке $x \in [a, b]$, имеющая вид

$$|f(x) - P_n(x)| \leq \frac{M_{n+1}}{(n+1)!} |\omega_{n+1}(x)|, \quad (11.26)$$

а также оценка максимума модуля погрешности интерполяции на отрезке $[a, b]$, имеющая вид

$$\max_{[a, b]} |f(x) - P_n(x)| \leq \frac{M_{n+1}}{(n+1)!} \max_{[a, b]} |\omega_{n+1}(x)|. \quad (11.27)$$

Здесь $M_{n+1} = \max_{[a, b]} |f^{(n+1)}(x)|$; предполагается, что производная $f^{(n+1)}$ непрерывна.

Пример 11.4. Оценим погрешность приближений к \ln (1.23), полученных в примере 11.3 с помощью интерполяции многочленами первой и второй степени. В этих случаях неравенство (11.26) примет вид

$$|f(x) - L_1(x)| \leq \frac{M_2}{2} |(x - x_0)(x - x_1)|, \quad (11.28)$$

$$|f(x) - L_2(x)| \leq \frac{M_3}{6} |(x - x_0)(x - x_1)(x - x_2)|. \quad (11.29)$$

Заметим, что для $f(x) = \ln(x)$ имеем $f^{(2)}(x) = -1/x^2$ и $f^{(3)}(x) = 2/x^3$. Поэтому здесь $M_2 = \max_{[1.2, 1.3]} |f^{(2)}(x)| = 1/1.2^2 \approx 0.69$ и $M_3 = \max_{[1.1, 1.3]} |f^{(3)}(x)| = 2/1.1^3 \approx 1.5$. Тогда в силу неравенств (11.28) и (11.29) получаем следующие оценки погрешности:

$$\epsilon_1 = |\ln(1.23) - L_1(1.23)| \lesssim \frac{0.69}{2} |(1.23 - 1.2)(1.23 - 1.3)| \approx 7.3 \cdot 10^{-4},$$

$$\begin{aligned} \epsilon_2 &= |\ln(1.23) - L_2(1.23)| \lesssim \frac{1.5}{6} |(1.23 - 1.1)(1.23 - 1.2)(1.23 - 1.3)| \approx \\ &\approx 6.9 \cdot 10^{-5}. \end{aligned}$$

Если на отрезке $[a, b]$ производная $f^{(n+1)}$ меняется слабо, то значение абсолютной погрешности $|f(x) - P_n(x)|$ почти полностью определяется значением функции $\omega_{n+1}(x)$. Представление о типичном характере поведения этой функции можно получить из рис. 11.2.

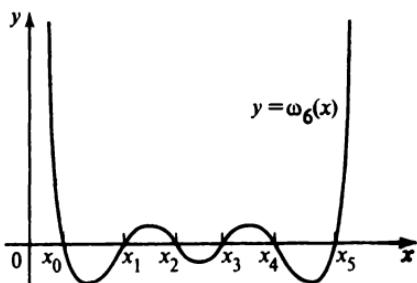


Рис. 11.2

Обратим внимание на то, что при выходе аргумента x за пределы отрезка наблюдения $[x_{\min}, x_{\max}]$ значение $|\omega_{n+1}(x)|$ быстро становится очень большим. Это объясняет ненадежность экстраполяции функции для значений аргумента, удаленных от отрезка наблюдения.

Пусть теперь $x_0 < x_1 < \dots < x_n$ и пусть $h_i = x_i - x_{i-1}$ — i -й шаг таблицы, а $h_{\max} = \max_{1 \leq i \leq n} h_i$. Несколько огрубив оценку (11.27), можно получить следующее неравенство:

$$\max_{[x_0, x_n]} |f(x) - P_n(x)| \leq \frac{M_{n+1}}{4(n+1)} h_{\max}^{n+1}. \quad (11.30)$$

Оно позволяет утверждать, что для достаточно гладкой функции f при фиксированной степени интерполяционного многочлена погрешность интерполяции на отрезке $[x_0, x_n]$ при $h_{\max} \rightarrow 0$ стремится к нулю не медленнее, чем некоторая величина, пропорциональная h_{\max}^{n+1} . Этот факт принято формулировать так: интерполяция многочленом степени n имеет $(n+1)$ -й порядок точности относительно h_{\max} . В частности, линейная и квадратичная интерполяции имеют второй и третий порядки точности соответственно.

§ 11.5. Интерполяция с кратными узлами

1. Интерполяционный многочлен с кратными узлами. Иногда в узлах x_i ($i = 0, 1, \dots, m$) бывают заданы не только значения $y_i = f(x_i)$ функции f , но и значения ее производных $y'_i = f'(x_i)$, $y''_i = f''(x_i)$, ..., $y^{(k_i-1)}_i = f^{(k_i-1)}(x_i)$ до некоторого порядка $k_i - 1$. В этом случае узел x_i называют *кратным*, а число k_i , равное количеству заданных

значений, — *кратностью узла*. Пусть $n = k_0 + k_1 + \dots + k_m - 1$. Можно доказать, что существует единственный многочлен $P_n(x)$ степени n , удовлетворяющий условиям

$$P_n(x_i) = y_i, \quad P'_n(x_i) = y'_i, \quad \dots, \quad P_n^{(k_i-1)}(x_i) = y_i^{(k_i-1)}$$

для всех $i = 0, 1, \dots, m$. Этот многочлен называют *интерполяционным многочленом с кратными узлами* или *интерполяционным многочленом Эрмита*¹. Можно указать и явную формулу его записи, аналогичную форме Лагранжа (11.22). Мы этого делать не будем и отметим лишь два важных частных случая.

1°. Пусть на концах отрезка $[x_0, x_1]$ заданы значения y_0, y_1, y'_0, y'_1 . Тогда $m = 1, k_0 = 2, k_1 = 2, n = 3$ и интерполяционный многочлен $P_3(x)$, удовлетворяющий условиям

$$P_3(x_0) = y_0, \quad P'_3(x_0) = y'_0, \quad P_3(x_1) = y_1, \quad P'_3(x_1) = y'_1,$$

может быть представлен (что проверяется непосредственно) в следующем виде:

$$\begin{aligned} P_3(x) = & y_0 \frac{(x_1 - x)^2(2(x - x_0) + h)}{h^3} + y'_0 \frac{(x_1 - x)^2(x - x_0)}{h^2} + \\ & + y_1 \frac{(x - x_0)^2(2(x_1 - x) + h)}{h^3} + y'_1 \frac{(x - x_0)^2(x - x_1)}{h^2}, \end{aligned} \quad (11.31)$$

где $h = x_1 - x_0$.

Многочлен (11.31) принято называть *кубическим интерполяционным многочленом Эрмита*.

2°. Пусть в точке x_0 заданы значения $y_0, y'_0, \dots, y_0^{(n)}$. Тогда многочлен $P_n(x)$, удовлетворяющий условиям $P_n(x_0) = y_0, P'_n(x_0) = y'_0, \dots, P_n^{(n)}(x_0) = y_0^{(n)}$ представляется в виде

$$P_n(x) = \sum_{k=0}^n y_0^{(k)} \frac{(x - x_0)^k}{k!}. \quad (11.32)$$

Как нетрудно видеть, многочлен $P_n(x)$ представляет собой отрезок ряда Тейлора. Таким образом, формула Тейлора дает решение задачи интерполяции² с одним узлом кратности $n + 1$.

¹ Шарль Эрмит (1822—1901) — французский математик.

² Заметим, что в действительности с ее помощью осуществляется экстраполяция.

2. Погрешность интерполяции с кратными узлами.

Теорема 11.5. Пусть функция f дифференцируема $n + 1$ раз на отрезке $[a, b]$, содержащем узлы интерполяции x_i ($i = 0, 1, \dots, m$). Тогда для погрешности интерполяции с кратными узлами в точке $x \in [a, b]$ справедливы равенство (11.25) и неравенства (11.26), (11.27), в которых $\omega_{n+1} = (x - x_0)^{k_0}(x - x_1)^{k_1} \dots (x - x_m)^{k_m}$, а ξ — некоторая точка, принадлежащая интервалу (a, b) . ■

Для формулы Тейлора ($m = 0$, $k_0 = n + 1$) теорема 11.5 дает известную формулу остаточного члена в форме Лагранжа. Для кубического многочлена Эрмита ($m = 1$, $k_0 = 2$, $k_1 = 2$) неравенство (11.27) приводит к следующей оценке погрешности:

$$\max_{[x_0, x_1]} |f(x) - P_3(x)| \leq \frac{M_4}{384} h^4. \quad (11.33)$$

Здесь учтено то, что максимум функции $\omega_4(x) = (x - x_0)^2(x - x_1)^2$ на отрезке $[x_0, x_1]$ достигается в точке $x = (x_0 + x_1)/2$ и равен $h^4/16$.

§ 11.6. Минимизация оценки погрешности интерполяции. Многочлены Чебышева

1. Постановка задачи минимизации оценки погрешности. Предположим, что значение заданной на отрезке $[a, b]$ функции f можно вычислить в произвольной точке x . Однако по некоторым причинам¹ целесообразнее заменить прямое вычисление функции f вычислением значений ее интерполяционного многочлена P_n . Для такой замены необходимо один раз получить таблицу значений функции f в выбранных на отрезке $[a, b]$ точках x_0, x_1, \dots, x_n . При этом естественно стремиться к такому выбору узлов интерполяции, который позволит сделать минимальной величину $\Delta(P_n) = \max_{[a, b]} |f(x) - P_n(x)|$ — погрешность интерполяции на отрезке $[a, b]$.

Пусть о функции f известно лишь то, что она непрерывно дифференцируема $n + 1$ раз на отрезке $[a, b]$. Тогда неравенство (11.27) дает верхнюю границу погрешности интерполяции:

$$\bar{\Delta}(P_n) = \frac{M_{n+1}}{(n+1)!} \max_{[a, b]} |\omega_{n+1}(x)|. \quad (11.34)$$

¹Например, вычисление значений $f(x)$ — трудоемкая операция.

Поставим теперь следующую задачу: определить набор узлов интерполяции x_0, x_1, \dots, x_n , при котором величина $\bar{\Delta}(P_n)$ минимальна. Для решения этой задачи нам потребуются некоторые сведения о многочленах Чебышёва¹.

Замечание. Формула (11.34) остается справедливой и в случае, когда некоторые из узлов x_0, x_1, \dots, x_n совпадают, т.е. имеет место интерполяция с кратными узлами.

2. Многочлены Чебышёва. Введенные П.Л. Чебышёвым многочлены $T_n(x)$ широко используются в вычислительной математике. При $n = 0$ и $n = 1$ они определяются явными формулами

$$T_0(x) = 1, \quad T_1(x) = x, \quad (11.35)$$

а при $n \geq 2$ рекуррентной формулой

$$T_n(x) = 2xT_{n-1}(x) - T_{n-2}(x). \quad (11.36)$$

Запишем явные формулы для многочленов Чебышёва $T_n(x)$ при $n = 2, 3, 4, 5$:

$$T_2(x) = 2xT_1(x) - T_0(x) = 2x^2 - 1,$$

$$T_3(x) = 2xT_2(x) - T_1(x) = 4x^3 - 3x,$$

$$T_4(x) = 2xT_3(x) - T_2(x) = 8x^4 - 8x^2 + 1,$$

$$T_5(x) = 2xT_4(x) - T_3(x) = 16x^5 - 20x^3 + 5x.$$

Аналогично можно записать явные формулы и при $n \geq 6$.

Приведем некоторые свойства многочленов Чебышёва.

1°. При четном n многочлен $T_n(x)$ содержит только четные степени x и является четной функцией, а при нечетном n многочлен $T_n(x)$ содержит только нечетные степени x и является нечетной функцией.

2°. При $n \geq 1$ старший коэффициент многочлена $T_n(x)$ равен 2^{n-1} , т.е. $T_n(x) = 2^{n-1}x^n + \dots$.

Справедливость свойств 1° и 2° следует непосредственно из определения (11.35), (11.36).

3°. Для $x \in [-1, 1]$ справедлива формула

$$T_n(x) = \cos(n \arccos x). \quad (11.37)$$

¹ Пафнутий Львович Чебышёв (1821—1894) — русский математик, один из создателей современных теории чисел, теории вероятностей, теории приближений функций.

Доказательство. При $n = 0$ и $n = 1$ формула (11.37) верна, так как $\cos(0 \cdot \arccos x) = 1$, $\cos(1 \cdot \arccos x) = x$. Для того чтобы доказать справедливость формулы для всех $n \geq 0$, достаточно показать, что функции $C_n(x) = \cos(n \arccos x)$ удовлетворяют такому же, как и многочлены Чебышева, рекуррентному соотношению

$$C_n(x) = 2x C_{n-1}(x) - C_{n-2}(x) \quad (11.38)$$

(сравните с (11.36)). Соотношение (11.38) получится, если в легко проверяемом тригонометрическом тождестве

$$\cos[(m+1)\varphi] + \cos[(m-1)\varphi] = 2\cos\varphi\cos m\varphi$$

положить $m = n - 1$ и $\varphi = \arccos x$. ■

4°. При $n \geq 1$ многочлен $T_n(x)$ имеет равно n действительных корней, расположенных на отрезке $[-1, 1]$ и вычисляемых по формуле

$$x_k = \cos \frac{(2k+1)\pi}{2n}, \quad k = 0, 1, \dots, n-1. \quad (11.39)$$

5°. При $n \geq 0$ справедливо равенство $\max_{[-1, 1]} |T_n(x)| = 1$. Если $n \geq 1$,

то этот максимум достигается ровно в $n+1$ точках, которые находятся по формуле

$$x_m = \cos \frac{\pi m}{n}, \quad m = 0, 1, \dots, n. \quad (11.40)$$

При этом $T_n(x_m) = (-1)^m$, т.е. максимумы и минимумы многочлена Чебышева чередуются.

Доказательство свойств 4° и 5° основано на применении формулы (11.37). Например, в силу этой формулы корни многочлена $T_n(x)$, расположенные на отрезке $[-1, 1]$, совпадают с корнями уравнения $\cos(n \arccos x) = 0$. Эквивалентное преобразование этого уравнения дает $n \cdot \arccos x = \pi/2 + \pi k$, $k = 0, \pm 1, \pm 2, \dots$. Так как $0 \leq \arccos x \leq \pi$, то заключаем, что имеется ровно n корней x_k , отвечающих значениям $k =$

$= 0, 1, \dots, n-1$ и удовлетворяющих равенствам $\arccos x_k = \frac{(2k+1)\pi}{2n}$, эквивалентным формуле (11.39).

Назовем величину $\max_{[-1, 1]} |P_n(x)|$ уклонением многочлена $P_n(x)$ от нуля. Эта величина характеризует максимальное отклонение (уклонение) графика многочлена P_n от графика функции $y = 0$ на отрезке $[-1, 1]$.

6°. Среди всех многочленов фиксированной степени $n \geq 1$ со старшим коэффициентом a_n , равным 1, наименьшее уклонение от нуля (равное 2^{1-n}) имеет многочлен $\bar{T}_n(x) = 2^{1-n}T_n(x)$.

Благодаря этому свойству, имеющему особую ценность для приложений, многочлены Чебышева иногда называют *наименее уклоняющимися от нуля*. Свойство 6° иначе можно сформулировать так: для любого многочлена вида $P_n(x) = x^n + a_{n-1}x^{n-1} + \dots + a_0$, отличного от $\bar{T}_n(x)$, справедливо неравенство

$$2^{1-n} = \max_{[-1, 1]} |\bar{T}_n(x)| < \max_{[-1, 1]} |P_n(x)|.$$

Приведем графики многочленов $T_n(x)$ для $n = 1, 2, 3, 4, 5$ (рис. 11.3).

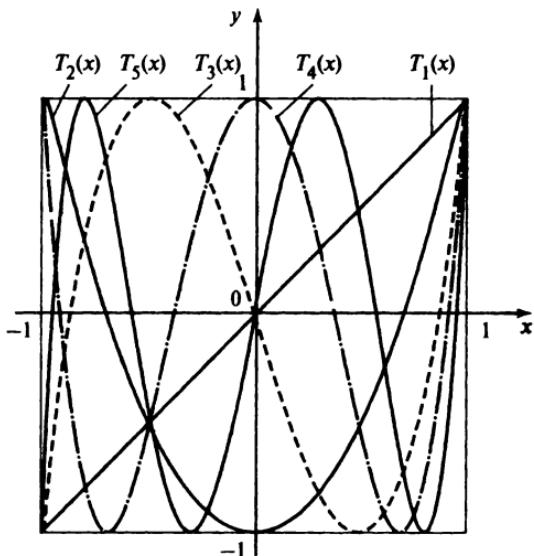


Рис. 11.3

Замечание Из свойства 6° следует, что среди всех многочленов $P_n(x)$ фиксированной степени $n \geq 1$ со старшим коэффициентом $a_n \neq 0$ наименьшее уклонение от нуля (равное $|a_n|2^{1-n}$) имеет многочлен $a_n\bar{T}_n(x)$.

Формулы (11.39) и (11.40) позволяют дать следующую геометрическую интерпретацию построения корней и точек экстремума много-

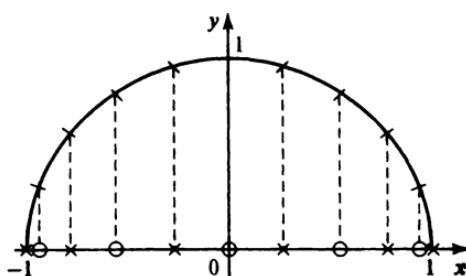


Рис. 11.4

члена $T_n(x)$. Разделим полуокружность, опирающуюся на отрезок $[-1, 1]$ как на диаметр, на $2n$ равных частей и спроектируем полученные точки на отрезок $[-1, 1]$. На рис. 11.4 изображен случай $n = 5$.

Нумеруя проекции справа налево, получим, что все проекции с нечетными номерами являются корнями многочлена T_n (на рис. 11.4 они помечены кружочками), а все проекции с четными номерами — точками экстремума (они помечены крестиками). Заметим, что корни и точки экстремума сгущаются к концам отрезка $[-1; 1]$.

3. Решение задачи минимизации оценки погрешности. Найдем сначала решение задачи в предположении, что отрезок интерполяции $[a, b]$ совпадает с отрезком $[-1, 1]$. В этом случае величина (11.34) будет минимальной при таком выборе узлов x_0, x_1, \dots, x_n , при котором минимальна величина $\max_{[-1, 1]} |\omega_{n+1}(x)|$, т.е. минимально уклонение многочлена $\omega_{n+1}(x) = (x - x_0)(x - x_1)\dots(x - x_n)$ от нуля. В силу свойств 4° и 6° многочленов Чебышева решение задачи дает набор узлов

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

являющихся нулями многочлена T_{n+1} , так как в этом случае $\omega_{n+1} = \bar{T}_{n+1}$.

Заметим, что при таком выборе

$$\bar{\Delta}(P_n) = \frac{M_{n+1}}{(n+1)!2^n}, \quad (11.41)$$

причем в силу свойства 6° любой другой выбор узлов дает большее значение верхней границы погрешности. Для сравнения укажем, что

при использовании для приближения функции f отрезка ряда Тейлора

$P_n(x) = \sum_{k=0}^n \frac{f^{(k)}(0)}{k!} x^k$ верхняя граница оценки погрешности такова:

$$\bar{\Delta}(P_n) = \frac{M_{n+1}}{(n+1)!}.$$

Следовательно, она в 2^n раз хуже, чем при интерполяции с оптимальным выбором узлов.

Пусть теперь отрезок интерполяции $[a, b]$ произволен. Приведем его к стандартному отрезку $[-1, 1]$ заменой

$$x = \frac{a+b}{2} + \frac{b-a}{2} t, \quad (11.42)$$

где $t \in [-1, 1]$. Как нетрудно видеть, в этом случае $\omega_{n+1}(x) = \left[\frac{b-a}{2} \right]^{n+1} \tilde{\omega}_{n+1}(t)$, где $\tilde{\omega}_{n+1}(t) = (t - t_0)(t - t_1) \dots (t - t_n)$ и $x_k = \frac{a+b}{2} + \frac{b-a}{2} t_k$ для $k = 0, 1, \dots, n$. Следовательно,

$$\bar{\Delta}(P_n) = \frac{M_{n+1}}{(n+1)!} \left[\frac{b-a}{2} \right]^{n+1} \max_{[-1, 1]} |\tilde{\omega}_{n+1}(t)|$$

и минимум этой величины достигается при значениях t_0, t_1, \dots, t_n , совпадающих с нулями многочлена T_{n+1} . Значит, решение поставленной задачи дает выбор узлов

$$x_k = \frac{a+b}{2} + \frac{b-a}{2} \cos \left(\frac{2k+1}{2n+2} \pi \right), \quad k = 0, 1, \dots, n, \quad (11.43)$$

которому отвечает минимальное значение верхней границы погрешности интерполяции, равное

$$\bar{\Delta}(P_n) = \frac{M_{n+1}}{(n+1)! 2^n} \left[\frac{b-a}{2} \right]^{n+1}.$$

§ 11.7. Конечные разности

1. Таблица конечных разностей. Пусть функция $y = f(x)$ задана таблицей (11.1) своих значений, причем $x_0 < x_1 < \dots < x_n$ и расстояние $h = x_i - x_{i-1}$ между соседними узлами таблицы значений аргумента постоянно. В этом случае величину h называют *шагом таблицы*, а узлы — *равноотстоящими*

Величину $\Delta y_i = y_{i+1} - y_i$ принято называть *конечной разностью первого порядка* функции $y = f(x)$ в точке x_i (с шагом h). *Конечная разность второго порядка* определяется формулой $\Delta^2 y_i = \Delta y_{i+1} - \Delta y_i$. Аналогично определяются конечные разности третьего и более высоких порядков. Общее определение *конечной разности порядка k* таково.

$$\Delta^k y_i = \Delta^{k-1} y_{i+1} - \Delta^{k-1} y_i;$$

здесь $k \geq 1$ и $\Delta^0 y_i = y_i$.

Таблицу *конечных разностей* (которые называют еще *конечными разностями вперед*) обычно располагают следующим образом (табл. 11.2):

Таблица 11.2

x	y	Δy	$\Delta^2 y$	$\Delta^3 y$		$\Delta^n y$
x_0	y_0					
x_1	y_1	Δy_0	$\Delta^2 y_0$	$\Delta^3 y_0$		
x_2	y_2	Δy_1	$\Delta^2 y_1$	$\Delta^3 y_1$		
x_3	y_3	Δy_2	$\Delta^2 y_2$	$\Delta^3 y_2$		
•	•		•	•		
•	•	•	•	•		
•	•	•	•	•		
						$\Delta^n y_0$
x_{n-1}	y_{n-1}	Δy_{n-2}	$\Delta^2 y_{n-2}$	$\Delta^3 y_{n-3}$		
x_n	y_n	Δy_{n-1}				

2. Свойства конечных разностей. Можно показать, что конечные разности порядка k выражаются через значения функции в $k + 1$ точке по формуле

$$\Delta^k y_i = \sum_{l=0}^k (-1)^{k-l} C_k^l y_{i+l}, \quad (11.44)$$

где $C_k^l = \frac{k!}{l!(k-l)!}$ — биномиальные коэффициенты. В частности,

$$\Delta^2 y_i = y_{i+2} - 2y_{i+1} + y_i,$$

$$\Delta^3 y_i = y_{i+3} - 3y_{i+2} + 3y_{i+1} - y_i,$$

$$\Delta^4 y_i = y_{i+4} - 4y_{i+3} + 6y_{i+2} - 4y_{i+1} + y_i.$$

Приведем без доказательства важное утверждение, указывающее на тесную связь между производными гладких функций и их конечными разностями.

Теорема 11.6. Пусть функция f дифференцируема k раз на отрезке $[x_i, x_{i+k}]$. Тогда справедливо равенство

$$\Delta^k y_i = h^k f^{(k)}(\xi), \quad (11.45)$$

в котором ξ — некоторая точка из интервала (x_i, x_{i+k}) . ■

Замечание. При $k=1$ формула (11.45) совпадает с формулой конечных приращений Лагранжа.

Следствие. Для многочлена $y = P_n(x) = \sum_{m=0}^n a_m x^m$ конечная разность порядка n является постоянной величиной, равной $h^n n! a_n$. Разности порядка $k > n$ тождественно равны нулю.

Конечные разности имеют разнообразные практические применения. Например, если производная k -го порядка $f^{(k)}$ слабо меняется на отрезке $[x_i, x_{i+k}]$, то в силу равенства (11.45) для $x \in [x_i, x_{i+k}]$ справедлива следующая формула численного дифференцирования:

$$f^{(k)}(x) \approx \frac{\Delta^k y_i}{h^k}. \quad (11.46)$$

В § 11.9 конечные разности будут использованы для построения интерполяционного многочлена Ньютона. Рассмотрим еще два приложения конечных разностей, связанных с анализом погрешностей таблиц, а именно задачу об оценке уровня «шума» таблицы и задачу обнаружения единичных ошибок.

Заметим, что в реальных вычислениях таблица конечных разностей $\Delta^k y_i$ строится по значениям y_j^* , каждое из которых содержит погреш-

ность $\varepsilon_j = y_j - y_j^*$. Тогда в силу формулы (11.44) найденные значения $\Delta^k y_j^*$ содержат неустранимые ошибки

$$\varepsilon_i^{(k)} = \Delta^k y_i - \Delta^k y_i^* = \sum_{l=0}^k (-1)^{k-l} C_k^l \varepsilon_{i+l}. \quad (11.47)$$

Как нетрудно видеть, имеется тенденция к росту ошибки $\varepsilon_i^{(k)}$ с ростом k . Если известно, что $|\varepsilon_i| \leq \varepsilon$ для всех i , то можно гарантировать справедливость лишь следующей оценки:

$$|\varepsilon_i^{(k)}| \leq \sum_{l=0}^k C_k^l \varepsilon = 2^k \varepsilon. \quad (11.48)$$

3. Оценка уровня «шума» в таблице. На практике часто возникает следующая задача. Для набора x_0, x_1, \dots, x_n равноотстоящих узлов каким-либо образом построена таблица приближенных значений гладкой функции $y = f(x)$. Требуется оценить уровень погрешности (уровень «шума») таблицы.

Полученная выше гарантированная оценка погрешности (11.48) не дает удовлетворительного ответа на поставленный вопрос. Она лишь указывает на то, что в самом неблагоприятном случае рост ошибки произойдет с коэффициентом, равным 2^k . Проведем статистический анализ погрешности. Будем предполагать, что ошибки ε_j ($j = 0, 1, \dots, n$) являются независимыми случайными величинами с математическим ожиданием $M[\varepsilon_j] = 0$ (это эквивалентно отсутствию систематической составляющей погрешности) и дисперсией $M[\varepsilon_j^2] = \sigma^2$.

В силу формулы (11.47) $M[\varepsilon_i^{(k)}] = 0$ и тогда для дисперсии $(\sigma^{(k)})^2 = M[(\varepsilon_i^{(k)})^2]$ погрешности k -й разности имеем:

$$\begin{aligned} (\sigma^{(k)})^2 &= M\left[\left(\sum_{l=0}^k (-1)^{k-l} C_k^l \varepsilon_{i+l}\right)\left(\sum_{r=0}^k (-1)^{k-r} C_k^r \varepsilon_{i+r}\right)\right] = \\ &= M\left[\sum_{l=0}^k (C_k^l)^2 \varepsilon_{i+l}^2 + \sum_{l \neq r} (-1)^{l+r} C_k^l C_k^r \varepsilon_{i+l} \varepsilon_{i+r}\right]. \end{aligned}$$

Заметим, что $\sum_{l=0}^k (C_k^l)^2 = C_{2k}^k$ и $M[\varepsilon_{i+l}\varepsilon_{i+r}] = 0$ при $l \neq r$, так как

величины ε_{i+l} и ε_{i+r} независимы. Следовательно, $(\sigma^{(k)})^2 = C_{2k}^k \sigma^2$.

Принимая за уровень «шума» таблицы величину σ квадратного корня из дисперсии (среднеквадратичную ошибку), получаем равенство

$$\sigma^{(k)} = \sqrt{C_{2k}^k} \sigma, \quad (11.49)$$

которое дает более оптимистичное по сравнению с оценкой (11.47) значение коэффициента роста ошибки, так как $\sqrt{C_{2k}^k} < 2^k$.

Если конечные разности $\Delta^k y_i^*$ строятся для гладкой функции, то они часто имеют тенденцию с ростом k уменьшаться по абсолютному значению, а затем, начиная с некоторого $k = p$, возрастать, испытывая сильные колебания в пределах одного столбца. При этом для $k \geq p$ основной вклад в значение $\Delta^k y_i^*$ вносит величина $\varepsilon_i^{(k)}$. Это обстоятельство позволяет считать, что в оценке дисперсии $(\sigma^{(p)})^2 \approx \frac{1}{n-p+1} \sum_{i=0}^{n-p} (\varepsilon_i^{(p)})^2$ величины $\varepsilon_i^{(p)}$ можно приближенно заменить на

$\Delta^p y_i^*$. Вычислив $\sigma^{(p)}$, затем достаточно воспользоваться формулой (11.49) для оценки уровня «шума» таблицы.

Пример 11.5. Оценим уровень «шума» в таблице значений функции $y = \ln x$, заданной с шагом $h = 0.1$ (первый и второй столбцы табл. 11.3). Составляя таблицу конечных разностей, замечаем, что, начиная с $k = 5$, абсолютные значения разностей $\Delta^k y_j$ начинают возрастать. Оценим $\sigma^{(5)}$ следующим образом:

$$\sigma^{(5)} \approx \left[\frac{1}{6} \sum_{i=0}^5 (\Delta^5 y_i)^2 \right]^{1/2} \approx 3.5 \cdot 10^{-4}.$$

Учитывая, что $C_{10}^5 = 252$, имеем $\sigma = \sigma^{(5)}/\sqrt{252} \approx 2 \cdot 10^{-5}$. Таким образом, погрешность таблицы составляет примерно две единицы пятого разряда, а шестой и седьмой разряды уже не содержат полезной информации. Если бы таблицу предполагалось использовать на практике, то, по-видимому, имело бы смысл округлить значения y_i до пяти значащих цифр после десятичной точки.

Таблица 11.3

x	y	Δy	$\Delta^2 y$	$\Delta^3 y$	$\Delta^4 y$	$\Delta^5 y$
1.0	0.0000000	0.0953274				
1.1	0.0953274	0.0869958	-0.0083316	0.0013874		
1.2	0.1823232	0.0800516	-0.0069442	0.0009812	-0.0004062	0.0003216
1.3	0.2623748	0.0740886	-0.0059630	0.0008966	-0.0000846	-0.0002525
1.4	0.3364634	0.0690222	-0.0050664	0.0005595	-0.0003371	0.0004161
1.5	0.4054856	0.0645153	-0.0045069	0.0006385	0.0000790	-0.0003636
1.6	0.4700009	0.0606469	-0.0038684	0.0003539	-0.0002846	0.0003992
1.7	0.5306478	0.0571324	-0.0035145	0.0004685	0.0001146	-0.0003473
1.8	0.5877802	0.0540864	-0.0030460	0.0002358	-0.0002327	
1.9	0.6418666	0.0512752	-0.0028102			
2.0	0.6931428					



4. Обнаружение единичных ошибок. Анализ таблиц конечных разностей позволяет в некоторых случаях обнаруживать грубые единичные ошибки в таблицах гладких функций и даже частично их устранять. Прежде чем продемонстрировать сказанное на примере рассмотрим, как распространяется в таблице конечных разностей ошибка ε , допущенная только в одном значении y_i . Пользуясь формулой (11.47) или равенством $\varepsilon_j^{(k)} = \varepsilon_{j+1}^{(k-1)} - \varepsilon_j^{(k-1)}$ (где $\varepsilon_j^{(0)} = \varepsilon$ при $j = i$ и $\varepsilon_j^{(0)} = 0$ при $j \neq i$), получаем следующую таблицу распространения единичной ошибки, т.е. ошибки, допущенной в одной точке (табл. 11.4).

Пример 11.6. Пусть на отрезке $[1.5, 2.8]$ задана таблица значений функции $y = \ln x$ (первый и второй столбцы табл. 11.5). Составим для

Таблица 11.4

x	$\varepsilon^{(0)}$	$\varepsilon^{(1)}$	$\varepsilon^{(2)}$	$\varepsilon^{(3)}$	$\varepsilon^{(4)}$...	$\varepsilon^{(k)}$
...	ε
...
x_{i-2}	0		0		ε	...	$-C_k^{k-1}\varepsilon$
		0		ε		...	
x_{i-1}	0		ε		-4ε	...	$C_k^{k-2}\varepsilon$
		ε		-3ε		...	
x_i	ε		-2ε		6ε	...	
		$-\varepsilon$		3ε		...	
x_{i+1}	0		ε		-4ε	...	$(-1)^{k-2}C_k^2\varepsilon$
		0		$-\varepsilon$...	
x_{i+2}	0		0		ε	...	$(-1)^{k-1}C_k^1\varepsilon$
...	
...	$(-1)^k\varepsilon$

нее таблицу конечных разностей. Аномальное поведение разностей третьего и четвертого порядка указывает на наличие в таблице значений функции ошибки. Сравнение с табл. 11.4 распространения единичной ошибки приводит к заключению о том, что погрешность допущена в значении y , отвечающем $x = 2.2$. Тогда погрешностям ε , -4ε , 6ε , -4ε , ε табл. 11.4 приближенно отвечают значения $262 \cdot 10^{-6}$, $-1231 \cdot 10^{-6}$, $1774 \cdot 10^{-6}$, $-1222 \cdot 10^{-6}$, $282 \cdot 10^{-6}$ табл. 11.5. Это соответствие имеет место при $\varepsilon \approx 298 \cdot 10^{-6}$. Следовательно, табличное значение 0.788757 нужно заменить на 0.788459. ▲

Замечание. Часто вместо конечных разностей вперед $\Delta^k y_i$ используют *конечные разности назад*, определяемые рекуррентной формулой

$$\nabla^k y_i = \nabla^{k-1} y_i - \nabla^{k-1} y_{i-1}.$$

Здесь $k \geq 1$, $\nabla^0 y_i = y_i$. Заметим, что разности вперед и назад связаны равенством

$$\Delta^k y_i = \nabla^k y_{i+k}$$

Таблица 11.5

x	y	Δy	$\Delta^2 y$	$\Delta^3 y$	$\Delta^4 y$
1.5	0.405465	0.064539			
1.6	0.470004	0.060624	-0.003915	0.000450	
1.7	0.530628	0.057159	-0.003465	0.000373	-0.000077
1.8	0.587787	0.054067	-0.003092	0.000318	-0.000055
1.9	0.641854	0.051293	-0.002774	0.000271	-0.000047
2.0	0.693147	0.048790	-0.002503	0.000533	0.000262
2.1	0.741937	0.046820	-0.001970	-0.000698	-0.001231
2.2	0.788757	0.044152	-0.002668	0.001076	0.001774
2.3	0.832909	0.042560	-0.001592	-0.000146	-0.001222
2.4	0.875469	0.040822	-0.001738	0.000136	0.000282
2.5	0.916291	0.039220	-0.001602	0.000123	-0.000013
2.6	0.955511	0.037741	-0.001479	0.000105	-0.000018
2.7	0.993252	0.036367	-0.001374		
2.8	1.029619				

§ 11.8. Разделенные разности

1. **Таблица разделенных разностей.** Пусть функция f задана на таблице x_0, x_1, \dots, x_n значений аргумента с произвольным (не обязательно постоянным) шагом, причем точки таблицы занумерованы в произвольном (не обязательно возрастающем) порядке. Величины

$$f(x_i; x_{i+1}) = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}$$

принято называть *разделенными разностями первого порядка* функции f . *Разделенные разности второго порядка* определяются формулой

$$f(x_i; x_{i+1}; x_{i+2}) = \frac{f(x_{i+1}; x_{i+2}) - f(x_i; x_{i+1})}{x_{i+2} - x_i}.$$

Аналогично определяются разделенные разности третьего и более высоких порядков. Общее определение *разделенной разности порядка $k \geq 2$* таково:

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

Таблицу разделенных разностей обычно располагают следующим образом (табл. 11.6).

Таблица 11.6

x	$f(x)$	Разделенные разности порядка			
		первого	второго	...	n -го
x_0	$f(x_0)$				
x_1	$f(x_1)$	$f(x_0; x_1)$	$f(x_0; x_1; x_2)$.	.
x_2	$f(x_2)$	$f(x_1; x_2)$.	.	$f(x_0; x_1; \dots; x_n)$
.
.	.	.	$f(x_{n-2}; x_{n-1}; x_n)$.	.
x_n	$f(x_n)$	$f(x_{n-1}; x_n)$			

2. Свойства разделенных разностей. Разделенные разности обладают рядом замечательных свойств. Перечислим без доказательства некоторые из них.

1°. *Разделенная разность $f(x_i; x_{i+1}; \dots; x_{i+k})$ является симметричной функцией своих аргументов $x_i, x_{i+1}, \dots, x_{i+k}$ (т.е. ее значение не меняется при любой их перестановке).*

2°. *Пусть функция f имеет на отрезке $[a, b]$, содержащем точки $x_i, x_{i+1}, \dots, x_{i+k}$, производную порядка k . Тогда справедливо равенство*

$$f(x_i; x_{i+1}; \dots; x_{i+k}) = \frac{f^{(k)}(\xi)}{k!}, \quad (11.50)$$

где ξ — некоторая точка, расположенная на интервале (a, b) .

3°. В случае, когда таблица значений аргумента имеет постоянный шаг h , разделенная и конечная разности связаны равенством

$$f(x_i; x_{i+1}; \dots; x_{i+k}) = \frac{\Delta^k y_i}{h^k k!}. \quad (11.51)$$

Пример 11.7. Приведем таблицу (табл. 11.7) разделенных разностей для функции, заданной табл. 11.1. Вычисления произведены на 6-разрядном десятичном компьютере.

Таблица 11.7

x	$f(x)$	Разделенные разности порядка			
		первого	второго	третьего	четвертого
1.0	0.000000				
1.1	0.095310	0.953100	-0.414900	0.221333	
1.2	0.182322	0.870120	<u>-0.348500</u>	<u>0.172667</u>	<u>-0.121665</u>
1.3	0.262364	<u>0.800420</u>	-0.296700		
1.4	0.336472	0.741080			

Перенумеруем теперь узлы, положив $x_0 = 1.2$, $x_1 = 1.3$, $x_2 = 1.1$, $x_3 = 1.4$, $x_4 = 1.0$. Тогда таблица разделенных разностей примет следующий вид (табл. 11.8).

Таблица 11.8

x	$f(x)$	Разделенные разности порядка			
		первого	второго	третьего	четвертого
1.2	0.182322	<u>0.800420</u>			
1.3	0.262364	0.835270	<u>-0.348500</u>	<u>0.172650</u>	
1.1	0.095310	0.803873	-0.313970	<u>0.197000</u>	<u>-0.121750</u>
1.4	0.336472	0.841180	-0.373070		
1.0	0.000000				

В табл. 11.8 подчеркнуты разделенные разности, которые совпадают (как и должно быть в силу свойства 1°) с точностью до вычислительной погрешности с соответствующими разделенными разностями из табл. 11.7 (они также подчеркнуты). ▲

§ 11.9. Интерполяционный многочлен Ньютона. Схема Эйткена

1. Интерполяционный многочлен Ньютона с разделенными разностями. Использовав разделенные разности, интерполяционный многочлен можно записать в следующем виде:

$$\begin{aligned} P_n(x) &= f(x_0) + f(x_0; x_1)(x - x_0) + f(x_0; x_1; x_2)(x - x_0)(x - x_1) + \dots \\ &\dots + f(x_0; x_1; \dots, x_n)(x - x_0)(x - x_1)\dots(x - x_{n-1}) = \\ &= \sum_{k=0}^n f(x_0; x_1; \dots, x_k) \omega_k(x). \end{aligned} \quad (11.52)$$

Здесь $\omega_0(x) \equiv 1$, $\omega_k(x) = (x - x_0)(x - x_1)\dots(x - x_{k-1})$ при $k \geq 1$.

Записанный в таком виде интерполяционный многочлен называют *интерполяционным многочленом Ньютона с разделенными разностями*.

Замечание 1. Отметим очевидную (с учетом равенства (11.50)) аналогию между формулой Ньютона (11.52) и формулой Тейлора (11.32).

Замечание 2. Формулу (11.25) для погрешности интерполяции в точке x , не являющейся узловой, можно уточнить следующим образом:

$$f(x) - P_n(x) = f(x_0; \dots; x_n; x) \omega_{n+1}(x). \quad (11.53)$$

Мы не приводим доказательства этой замечательной формулы. Отметим лишь, что если воспользоваться свойством 2° разделенных разностей, то из нее немедленно получается формула (11.25).

В практическом плане формула (11.52) обладает рядом преимуществ перед формулой Лагранжа. Пусть, например, по каким-либо причинам необходимо увеличить степень интерполяционного многочлена на единицу, добавив в таблицу еще один узел x_{n+1} . При использовании формулы Лагранжа (11.22) это приводит не только к увеличению числа слагаемых, но и к необходимости вычислять каждое из них заново.

В то же время для вычисления $P_{n+1}(x)$ по формуле Ньютона (11.52) достаточно добавить к $P_n(x)$ лишь одно очередное слагаемое, так как

$$P_{n+1}(x) - P_n(x) = f(x_0; \dots; x_n; x_{n+1})\omega_{n+1}(x). \quad (11.54)$$

Заметим, что когда величина $|x_{n+1} - x|$ мала, а функция f достаточно гладкая, справедливо приближенное равенство

$$f(x_0; \dots; x_n; x) \approx f(x_0; \dots; x_n; x_{n+1}),$$

из которого с учетом равенств (11.53) и (11.54) следует, что

$$f(x) - P_n(x) \approx P_{n+1}(x) - P_n(x).$$

Таким образом, величину

$$\varepsilon_n = |P_{n+1}(x) - P_n(x)| \quad (11.55)$$

можно использовать для практической оценки погрешности интерполяции.

Пример 11.8. По табл. 11.1 значений функции $y = \ln x$ из примера 11.3 найдем приближенное значение $\ln x$ при $x = 1.23$, используя интерполяционные многочлены Ньютона с разделенными разностями $P_k(x)$ для $k = 0, 1, \dots, 4$. Оценим при $k = 0, 1, 2, 3$ погрешность интерполяции по формуле (11.55).

Занумеруем узлы таблицы в следующем порядке: $x_0 = 1.2$, $x_1 = 1.3$, $x_2 = 1.1$, $x_3 = 1.4$, $x_4 = 1.0$, т.е. в порядке возрастания расстояния до точки $x = 1.23$. Соответствующие этой нумерации разделенные разности приведены в табл. 11.8 (мы используем только подчеркнутые разности).

Вычисления на 6-разрядном десятичном компьютере дают следующие значения:

$$\begin{aligned} P_0(x) &= 0.182322, \quad P_1(x) = P_0(x) + f(x_0; x_1)\omega_1(x) \approx 0.182322 + \\ &+ 0.80042 \cdot (1.23 - 1.2) \approx 0.206335, \quad \varepsilon_0 = |P_1(x) - P_0(x)| \approx 2.4 \cdot 10^{-2}; \\ P_2(x) &= P_1(x) + f(x_0; x_1; x_2)\omega_2(x) \approx 0.206335 - 0.3485(1.23 - 1.2)(1.23 - 1.3) \approx \\ &\approx 0.207067, \quad \varepsilon_1 = |P_2(x) - P_1(x)| \approx 7.3 \cdot 10^{-4}. \end{aligned}$$

Аналогично получаются значения $P_3(x) \approx 0.207020$; $\varepsilon_2 \approx 4.7 \cdot 10^{-5}$; $P_4(x) \approx 0.207014$, $\varepsilon_3 \approx 6 \cdot 10^{-6}$.

Если бы задача состояла в определении значения $\ln(1.23)$ с точностью $\varepsilon = 10^{-4}$, то вычисления следовало бы окончить после получения $\varepsilon_2 < \varepsilon$. Результат был бы таким: $\ln(1.23) \approx P_2(x) \approx 0.2071$. ▲

2. Интерполяция с использованием схемы Эйткена. Рассмотрим один из алгоритмов решения задачи интерполяции. Предполагается,

что задана таблица значений функции f . Требуется при заданном x вычислить с помощью интерполяции значение $f(x)$ с заданной точностью ε либо с максимально возможной при имеющейся информации точностью. Считается, что функция f достаточно гладкая.

Обозначим через $P_{(k, k+1, \dots, m)}(x)$ интерполяционный многочлен степени $m - k$ с узлами интерполяции x_k, x_{k+1}, \dots, x_m . В частности, положим $P_{(k)}(x) = y_k$. В этих обозначениях справедливо равенство

$$P_{(k, k+1, \dots, m+1)}(x) = \frac{P_{(k+1, \dots, m+1)}(x)(x-x_k) - P_{(k, \dots, m)}(x)(x-x_{m+1})}{x_{m+1} - x_k}. \quad (11.56)$$

Действительно, правая часть представляет собой многочлен степени $m + 1 - k$. Непосредственная проверка показывает, что этот многочлен совпадает с y_i в точках $x = x_i$ для $i = k, k+1, \dots, m+1$ и, значит, по определению равен $P_{(k, \dots, m+1)}(x)$.

Удобный и экономичный способ вычисления значения многочлена $P_n(x) = P_{(0, 1, \dots, n)}(x)$, лежащий в основе рассматриваемого алгоритма, дает *схема Эйткена*¹. Она заключается в последовательном вычислении с помощью формулы (11.56) элементов следующей таблицы (табл. 11.9).

Таблица 11.9

Значения интерполяционных многочленов, вычисляемых по схеме Эйткена, порядка				
нулевого	первого	второго		n -го
$P_{(0)}(x) = y_0$				
	$P_{(0, 1)}(x)$			
$P_{(1)}(x) = y_1$		$P_{(0, 1, 2)}(x)$		
			\ddots	
$P_{(2)}(x) = y_2$	$P_{(1, 2)}(x)$			
.	.	.		
.	.	.		
.	.	$P_{(n-2, n-1, n)}(x)$		
			\ddots	
$P_{(n)}(x) = y_n$	$P_{(n-1, n)}(x)$			

Для решения поставленной задачи интерполяции при заданном значении x узлы нумеруют в порядке возрастания их расстояния $|x - x_k|$ до точки x . Затем последовательно вычисляют значения $P_1(x), \varepsilon_0, P_2(x), \varepsilon_1, \dots, P_{m+1}(x), \varepsilon_m, \dots$. Если при некотором m оказывается, что $\varepsilon_m \leq \varepsilon$, то

¹ Александр Крэг Эйткен (1895—1967) — английский математик.

полагают $f(x) \approx P_m(x)$. Если же $\varepsilon_m > \varepsilon$ для всех m , то полагают $f(x) \approx P_k(x)$, где k — степень, при которой достигается минимум оценки погрешности: $\varepsilon_k = \min_{m \geq 0} \varepsilon_m$.

Пример 11.9. Для решения задачи из примера 11.8 воспользуемся схемой Эйткена. Здесь (как и в примере 11.8) $x_0 = 1.2$, $x_1 = 1.3$, $x_2 = 1.1$, $x_3 = 1.4$, $x_4 = 1.0$. После завершения вычислений табл. 11.9 принимает следующий вид (табл. 11.10).

Таблица 11.10

Значения интерполяционных многочленов, вычисляемых по схеме Эйткена, порядка			
нулевого	первого	второго	третьего
<u>0.182322</u>			
0.262364	<u>0.206335</u>		
0.095310	0.203895	<u>0.207067</u>	<u>0.207020</u>
0.336472	0.199814	0.206752	

Подчеркнутые числа дают те же, что и в примере 11.8, значения $P_k(x)$, $k = 0, 1, 2, 3$. Естественно, что теми же окажутся и значения ε_k . ▲

3. Интерполяционный многочлен Ньютона с конечными разностями. Пусть интерполируемая функция задана на таблице с постоянным шагом h (т.е. $x_i = x_0 + ih$, $i = 0, 1, \dots, n$). В этом случае, используя формулу (11.51) связи между разделенными и конечными разностями и вводя безразмерную переменную $t = (x - x_0)/h$, многочлен Ньютона (11.52) можно записать в следующем виде:

$$\begin{aligned} P_n(x) = P_n(x_0 + ht) &= y_0 + \frac{\Delta y_0}{1!} t + \frac{\Delta^2 y_0}{2!} t(t-1) + \\ &+ \frac{\Delta^3 y_0}{3!} t(t-1)(t-2) + \dots + \frac{\Delta^n y_0}{n!} t(t-1)\dots(t-n+1). \end{aligned} \quad (11.57)$$

Многочлен (11.57) называется *интерполяционным многочленом Ньютона с конечными разностями для интерполяции вперед*.

Заметим, что в формуле (11.57) используются только конечные разности, расположенные в верхней косой строке табл. 11.2. Можно использовать конечные разности, расположенные и в нижней косой

строке табл. 11.2, записав многочлен в виде *интерполяционного многочлена Ньютона с конечными разностями для интерполяции назад*:

$$\begin{aligned} P_n(x) = P_n(x_n + hq) &= y_n + \frac{\nabla y_n}{1!} q + \frac{\nabla^2 y_n}{2!} q(q+1) + \\ &+ \frac{\nabla^3 y_n}{3!} q(q+1)(q+2) + \dots + \frac{\nabla^n y_n}{n!} q(q+1)\dots(q+n-1), \quad (11.58) \end{aligned}$$

где $q = (x - x_n)/h$.

§ 11.10. Обсуждение глобальной полиномиальной интерполяции.

Понятие о кусочно-полиномиальной интерполяции

Пусть функция интерполируется на отрезке $[a, b]$. Метод решения этой задачи с помощью интерполяции единственным для всего отрезка многочленом $P_n(x)$ называют *глобальной полиномиальной интерполяцией*. При первом знакомстве с интерполяцией этот подход кажется привлекательным. В самом деле, неплохо иметь один многочлен, пригодный для приближения функции f во всех точках $x \in [a, b]$. В то же время известные результаты теории аппроксимации позволяют надеяться на то, что удастся приблизить функцию с любой требуемой точностью ε с помощью соответствующего выбора степени многочлена и узлов интерполяции на отрезке $[a, b]$. Приведем один такой классический результат.

Теорема 11.6 (аппроксимационная теорема Вейерштрасса¹). *Пусть функция f непрерывна на отрезке $[a, b]$. Тогда для любого $\varepsilon > 0$ существует полином $P_n(x)$ степени $n = n(\varepsilon)$ такой, что*

$$\max_{[a,b]} |f(x) - P_n(x)| < \varepsilon. \blacksquare$$

Замечание. В формулировке теоремы Вейерштрасса нет информации о том, как построить многочлен P_n . Данное в 1912 г. С.Н. Бернштейном² изящное доказательство этой теоремы дает явную формулу для приближающего многочлена.

¹ Карл Теодор Вильгельм Вейерштрасс (1815—1897) — немецкий математик, один из основоположников современного математического анализа и теории аналитических функций.

² Сергей Натаевич Бернштейн (1880—1968) — советский математик. Родился в Одессе, получил образование в Париже, работал в Харькове, Ленинграде и Москве. Получил ряд фундаментальных результатов по теории дифференциальных уравнений, теории приближения функций и теории вероятностей.

$$P_n(x) = \sum_{k=0}^n f(x_k) C_n^k t^k (1-t)^{n-k}, \quad t = \frac{x-a}{b-a},$$

использующую значения функции f в точках $x_k = a + (b-a) \frac{k}{n}$.

К сожалению, этот способ приближения приводит к многочленам очень высокой степени и для практического использования, как правило, не пригоден. Отметим также, что многочлен Бернштейна не является интерполяционным.

Существуют весьма веские причины, по которым глобальная интерполяция многочленами высокой степени в вычислительной практике, как правило, не применяется. Обсудим некоторые из этих причин.

1. Сходимость при увеличении числа узлов. Всегда ли можно добиться повышения точности интерполяции благодаря увеличению числа узлов (и соответственно степени n интерполяционного многочлена)? Хотя положительный ответ на этот вопрос напрашивается сам собой, не будем торопиться с выводами.

Уточним постановку задачи. Для того чтобы реализовать процесс интерполяции функции f многочленами возрастающей степени n , необходимо указать стратегию выбора при каждом n набора узлов интерполяции $x_0^{(n)}, x_1^{(n)}, \dots, x_n^{(n)}$. Такая стратегия задается указанием *интерполяционного массива* — треугольной таблицы

$$\begin{array}{ccccccc} & & x_0^{(0)} & & & & \\ & & x_0^{(1)} & x_1^{(1)} & & & \\ & & x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & & \\ & & \cdot & \cdot & \cdot & \cdot & \cdot \\ & & x_0^{(n)} & x_1^{(n)} & x_2^{(n)} & \dots & x_n^{(n)}, \end{array}$$

в каждой строке которой все $x_i^{(n)}$ различны и $x_i^{(n)} \in [a, b]$. Будем говорить, что при заданной стратегии выбора узлов *метод интерполяции сходится*, если $\max_{[a, b]} |f(x) - P_n(x)| \rightarrow 0$ при $n \rightarrow \infty$.

Рассмотрим сначала простейшую стратегию, состоящую в равномерном распределении на отрезке $[a, b]$ узлов интерполяции, т.е. в выборе $x_i^{(n)} = a + ih$ ($i = 0, 1, \dots, n$), где $h = (b-a)/n$. Следующий

пример показывает, что такая стратегия не может обеспечить сходимость интерполяции даже для очень гладких функций.

Пример 11.10 (пример Рунге¹). Используем глобальную полиномиальную интерполяцию с равномерным распределением узлов для приближения на отрезке $[-1, 1]$ следующей функции:

$$f(x) = \frac{1}{1 + 25x^2}. \quad (11.59)$$

Вычисления показывают, что при больших n интерполяция дает превосходные результаты в центральной части отрезка. В то же время вопреки ожиданиям последовательность $P_n(x)$ расходится при $n \rightarrow \infty$ для $0.73 \leq |x| \leq 1$. Соответствующая иллюстрация приведена на рис. 11.5. ▲

Равномерное распределение узлов интерполяции для функции Рунге (11.59) оказалось неудачным. Однако проблема сходимости для этой функции исчезает, если в качестве узлов интерполяции брать корни многочлена Чебышева $T_{n+1}(x)$. Существует ли единая для всех непрерывных на отрезке $[a, b]$ функций f стратегия выбора узлов интерполяции, гарантирующая ее сходимость? Отрицательный ответ на этот вопрос дает следующая теорема.

Теорема 11.7 (теорема Фабера²). *Какова бы ни была стратегия выбора узлов интерполяции, найдется непрерывная на отрезке $[a, b]$ функция f , для которой $\max_{[a, b]} |f(x) - P_n(x)| \rightarrow \infty$ при $n \rightarrow \infty$.* ■

Теорема Фабера отрицает существование единой для всех непрерывных функций стратегии выбора узлов интерполяции. Однако для гладких функций (а именно такие функции чаще всего и интерполируются) такая стратегия существует, о чем говорит следующая теорема.

Теорема 11.8. *Пусть в качестве узлов интерполяции на отрезке $[a, b]$ выбираются чебышевские узлы (11.43). Тогда для любой непрерывно дифференцируемой на отрезке $[a, b]$ функции f метод интерполяции сходится.* ■

¹ Карл Давид Тольме Рунге (1856—1927) — немецкий физик и математик. Историческую справку о его жизни и деятельности см. в [54].

² Жорж Фабер (1877—1966) — швейцарский математик.

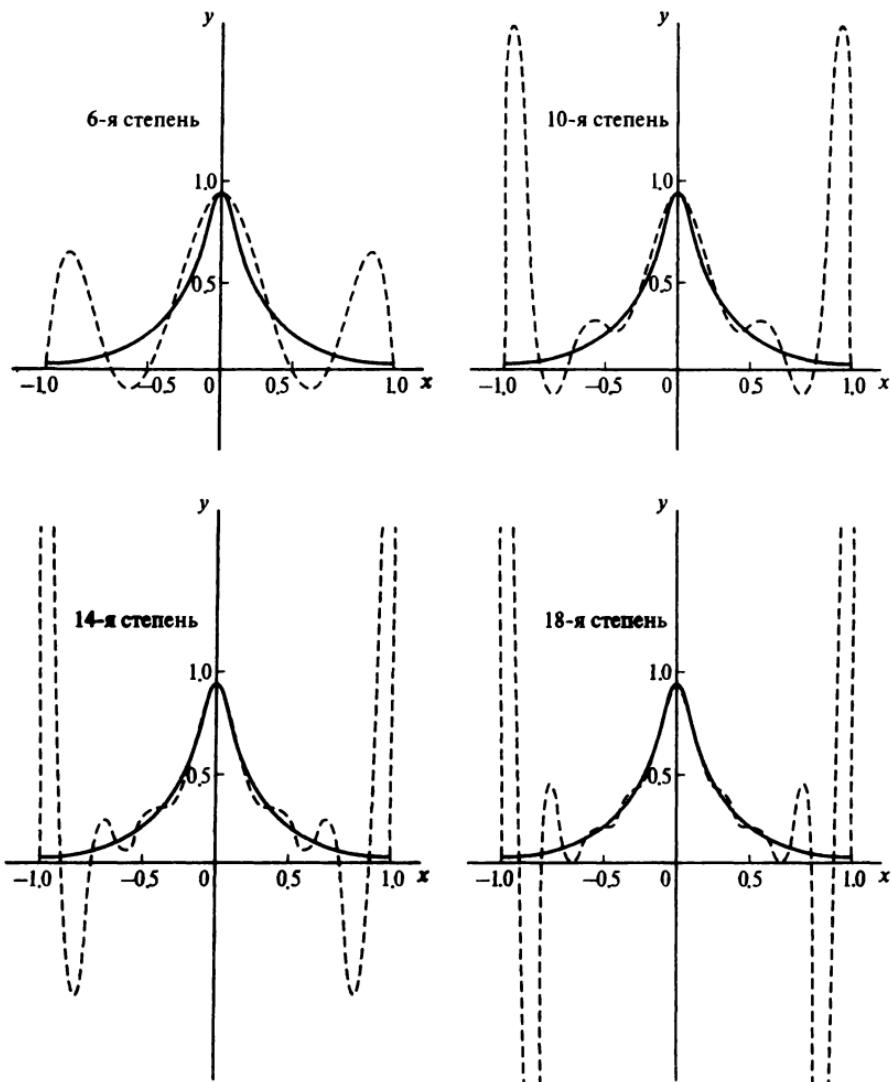


Рис. 11.5

Замечание. Практическая реализация стратегии выбора узлов интерполяции (11.43) возможна и оправдана в довольно редких случаях и просто невозможна тогда, когда приходится иметь дело с заданной таблицей значений функции.

2. Чувствительность интерполяционного многочлена к погрешностям входных данных. Помимо погрешности, которая возникает от приближенной замены функции f интерполяционным многочленом, возникает еще дополнительная погрешность, связанная с тем, что значения интерполируемой функции также задаются с погрешностью.

Пусть заданные в узлах x_i значения y_i^* содержат погрешности ε_i .

Тогда вычисляемый по этим значениям многочлен $P_n^*(x) = \sum_{j=0}^n y_j^* l_{nj}(x)$

содержит погрешность

$$P_n(x) - P_n^*(x) = \sum_{j=0}^n \varepsilon_j l_{nj}(x). \quad (11.60)$$

Например, при линейной интерполяции по приближенно заданным значениям справедливо равенство

$$P_1(x) - P_1^*(x) = \varepsilon_0 l_{10}(x) + \varepsilon_1 l_{11}(x),$$

где

$$l_{10} = \frac{x - x_1}{x_0 - x_1}, \quad l_{11}(x) = \frac{x - x_0}{x_1 - x_0}.$$

Воспользуемся тем, что $|l_{10}(x)| + |l_{11}(x)| = 1$ для $x \in [x_0, x_1]$. Следовательно,

$$\max_{[x_0, x_1]} |P_1(x) - P_1^*(x)| \leq \max\{|\varepsilon_0|, |\varepsilon_1|\}.$$

Таким образом, при линейной интерполяции погрешность, возникающая вследствие погрешности значений функции, не превосходит верхней границы погрешности этих значений.

Рассмотрим общий случай. Пусть известно, что верхняя граница погрешности значений y_i^* равна $\bar{\Delta}(y^*)$, т.е. $|\varepsilon_i| \leq \bar{\Delta}(y^*)$ для всех $i = 0, 1, \dots, n$. Тогда для верхней границы соответствующей погрешности интерполяционного многочлена $\bar{\Delta}(P_n^*) = \max_{[a, b]} |P_n(x) - P_n^*(x)|$ в силу равенства (11.60) справедлива оценка

$$\bar{\Delta}(P_n^*) \leq \Lambda_n \bar{\Delta}(y^*). \quad (11.61)$$

Здесь $\Lambda_n = \max_{[a, b]} \sum_{j=0}^n |l_{nj}(x)|$ — величина, которую называют *константой Лебега*¹.

Замечание. Если $|\varepsilon_i| = \bar{\Delta}(y^*)$ для всех i , то с помощью выбора знаков погрешностей ε_i можно добиться выполнения равенства $\bar{\Delta}(P_n^*) = \Lambda_n \bar{\Delta}(y^*)$. Это означает, что в самом неблагоприятном случае погрешность входных данных при интерполяции может возрасти в Λ_n раз. Таким образом, в задаче интерполирования константа Лебега играет роль абсолютного числа обусловленности.

Величина Λ_n не зависит от длины отрезка $[a, b]$, а определяется лишь относительным расположением узлов на нем. Для того чтобы показать это, приведем отрезок $[a, b]$ к стандартному отрезку $[-1, 1]$ с помощью линейного преобразования $x = \frac{a+b}{2} + \frac{b-a}{2} t$. Тогда $x_i = \frac{a+b}{2} + \frac{b-a}{2} t_i$ и константа Лебега приводится к виду

$$\Lambda_n = \max_{[-1, 1]} \sum_{j=0}^n |\bar{l}_{nj}(t)|, \text{ где } \bar{l}_{nj}(t) = \prod_{k \neq j} \frac{(t - t_k)}{(t_j - t_k)}.$$

Естественно поставить задачу о таком оптимальном выборе узлов интерполяции, при котором величина Λ_n оказалась минимально возможной. Для малых значений n эту задачу решить нетрудно. При $n = 1$ выбор $t_0 = -1, t_1 = 1$ дает $\Lambda_1 = 1$. При $n = 2$ оптимальное значение $\Lambda_2 = 5/4$ достигается, например, при $t_0 = -1, t_1 = 0, t_2 = 1$. В общем случае оптимальный выбор узлов неизвестен. Установлено, однако, что почти оптимальным является выбор в качестве узлов интерполяции нулей многочлена Чебышева T_{n+1} . При таком выборе $\Lambda_n \approx \frac{2}{\pi} \ln(n+1) + 1$.

С рассматриваемой точки зрения крайне неудачным при больших n является выбор равноотстоящих узлов интерполяции. При таком выборе $\Lambda_n > \frac{2^{n-1}}{(2n-1)\sqrt{n}}$ для $n \geq 4$ и обусловленность задачи резко ухудшается.

¹ Анри Леон Лебег (1875—1941) — французский математик, один из создателей современной теории функций вещественной переменной и теории интегрирования.

шается с ростом n . Сказанное позволяет сделать важный вывод: в вычислениях не следует использовать интерполяционные многочлены высокой степени с равноотстоящими узлами.

3. Обусловленность задачи вычисления многочлена, с приближенно заданными коэффициентами. Обратим внимание на еще один потенциальный источник потери точности при использовании многочленов $P_n(x)$ высокой степени. Для определенности будем считать, что многочлен вычисляется на отрезке $[a, b]$, где $|a| \leq b$, причем предварительно он представлен в виде

$$P_n(x) = \sum_{k=0}^n a_k \varphi_k(x). \quad (11.62)$$

Здесь $\{\varphi_k(x)\}_{k=0}^n$ — некоторый набор базисных многочленов $\varphi_k(x)$, обладающий тем свойством, что всякий многочлен $P_n(x)$ степени n может быть однозначно представлен в виде (11.62). Например, можно использовать степенной базис $\{x^k\}_{k=0}^n$, нормированный степенной базис

$\left\{ \left(\frac{x}{b} \right)^k \right\}_{k=0}^n$, локальный степенной базис $\left\{ \left(\frac{x-a}{b-a} \right)^k \right\}_{k=0}^n$, чебышевский базис

$\left\{ T_k \left(\frac{x-(a+b)/2}{(b-a)/2} \right) \right\}_{k=0}^n$, лагранжев базис $\{l_{nk}(x)\}_{k=0}^n$ (см. § 11.3) и т.д.

При вычислении коэффициентов a_k неизбежны погрешности, приводящие к приближенным значениям a_k^* . Поэтому в действительности будет вычисляться многочлен

$$P_n^*(x) = \sum_{k=0}^n a_k^* \varphi_k(x).$$

Примем за относительную погрешность вектора $\mathbf{a}^* = (a_0^*, a_1^*, \dots, a_n^*)^T$ величину

$$\delta(\mathbf{a}^*) = \max_{0 \leq i \leq n} |a_i - a_i^*| / \max_{0 \leq i \leq n} |a_i|,$$

а за относительную погрешность многочлена $P_n^*(x)$ — величину

$$\delta(P_n^*) = \max_{[a, b]} |P_n(x) - P_n^*(x)| / \max_{[a, b]} |P_n(x)|.$$

Числом обусловленности задачи вычисления многочлена с приближенно заданными коэффициентами назовем величину cond_n , равную минимальной из постоянных K_n , для которых выполняется неравенство

$$\delta(P_n^*) \leq K_n \delta(a^*).$$

Величина cond_n характеризует чувствительность вычисляемых значений многочлена к погрешностям в коэффициентах a_k и существенно зависит от выбора базиса $\{\phi_k(x)\}_{k=0}^n$. Неудачный выбор базиса может сделать эту, казалось бы, элементарную задачу очень плохо обусловленной. Доказано, например, что даже для вполне «разумного»

локального степенного базиса $\left\{ \left(\frac{x-a}{b-a} \right)^k \right\}_{k=0}^n$ число обусловленности

оценивается снизу величиной $T_n(3)$. Как видно из табл. 11.11, использование указанного базиса для представления многочленов высокой степени сопряжено с опасностью значительной потери точности.

Таблица 11.11

n	1	2	3	4	5	6	7	8	9
$T_n(3)$	3	17	99	577	3363	19601	114243	665857	3880899

Отметим, что для чебышевского базиса $\text{cond}_n \leq \sqrt{2}(n+1)$. Для лагранжева базиса число обусловленности совпадает с константой Лебега Λ_n . Напомним, что при почти оптимальном выборе узлов интерполяции $\Lambda_n \approx \frac{2}{\pi} \ln(n+1) + 1$.

Итак, глобальная полиномиальная интерполяция многочленом высокой степени может привести к неудаче или оказаться неэффективной. Альтернативный подход состоит в локальной интерполяции, когда функция f аппроксимируется интерполяционным многочленом $P_n(x)$ невысокой степени m на содержащемся в $[a, b]$ отрезке $[\alpha, \beta]$ малой длины. Естественно, что при этом используется лишь часть табличных значений. Рассмотрим два подхода к приближению функции, основанные на локальной интерполяции.

4. Интерполирование с помощью «движущегося» полинома. Строят набор полиномов $P_{(0, 1, \dots, m)}, P_{(1, 2, \dots, m+1)}, \dots, P_{(n-m, n-m+1, \dots, n)}$ фиксированной степени m , каждый из которых совпадает с табличными значениями в $m+1$ последовательных точках. Каждый такой полином используют для приближения функции в тех точках x из отрезка $[a, b]$, для которых выбранные узлы таблицы являются ближайшими.

Пример 11.11. Пусть функция задана следующей таблицей (табл. 11.12):

Таблица 11.12

i	0	1	2	3	4
x_i	0	1	2	3	4
v_i	1.0	1.8	2.2	1.4	1.0

Для интерполяции этой функции воспользуемся «движущимся» полиномом второй степени. Заметим, что при $x \in [0.0, 1.5]$ для приближения используется многочлен $P_{(0, 1, 2)}(x)$, при $x \in [1.5, 2.5]$ — многочлен $P_{(1, 2, 3)}(x)$, при $x \in [2.5, 4.0]$ — многочлен $P_{(2, 3, 4)}(x)$. Соответствующая геометрическая иллюстрация приведена на рис. 11.6. Заметим, что полученная таким способом аппроксимирующая функция имеет разрывы в точках $x = 1.5$ и $x = 2.5$. ▲

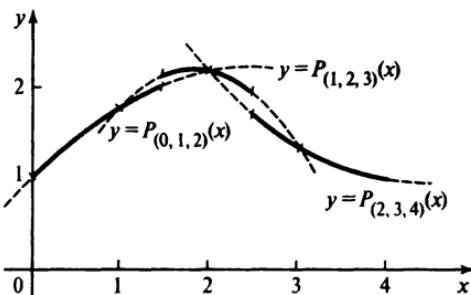


Рис. 11.6

5. Кусочно-полиномиальная интерполяция. Исходный отрезок $[a, b]$ разбивают на несколько отрезков меньшей длины, на каждом из которых функция интерполируется своим многочленом.

Пример 11.12. Для интерполяции функции из примера 11.11 используем кусочно-полиномиальную интерполяцию. На отрезке $[0, 2]$ аппроксимируем функцию многочленом $P_{(0, 1, 2)}(x)$, а на отрезке $[2, 4]$ — многочленом $P_{(2, 3, 4)}(x)$. Соответствующая геометрическая иллюстрация приведена на рис. 11.7. Заметим, что результирующая аппроксимирующая функция непрерывна, но в точке $x = 2$ график ее имеет излом, соответствующий разрыву первой производной. ▲

Заметим, что интерполяцию «движущимся» полиномом можно рассматривать как частный случай кусочно-полиномиальной интерполяции.

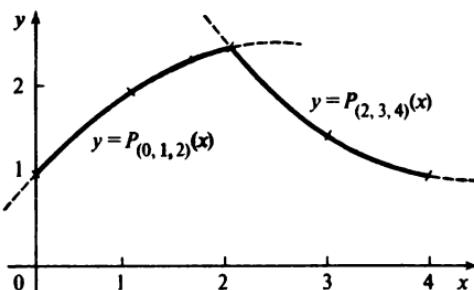


Рис. 11.7

Как следует из оценки (11.30), метод кусочно-полиномиальной интерполяции при использовании многочленов фиксированной степени m имеет $(m + 1)$ -й порядок точности относительно h_{\max} .

§ 11.11. Интерполяция сплайнами

1. Определение сплайна. Проведенное выше обсуждение интерполяции показывает, что повышение точности приближения гладкой функции благодаря увеличению степени интерполяционного многочлена возможно (см. теорему 11.8), но связано с существенным повышением сложности вычислений. К тому же использование многочленов высокой степени требует специальных мер предосторожности уже при выборе формы их записи, и вычисления сопровождаются накоплением погрешностей округления. Поэтому на практике предпочитают кусочно-полиномиальную интерполяцию с использованием многочленов невысокой степени. Однако этот способ приближения имеет недостаток: в точках «стыка» двух соседних многочленов производная, как правило, имеет разрыв (см. пример 11.12). Часто это обстоятельство не играет существенной роли. Вместе с тем нередко требуется, чтобы аппроксимирующая функция была гладкой и тогда простейшая кусочно-полиномиальная интерполяция становится неприемлемой.

Естественная потребность в наличии аппроксимирующих функций, которые сочетали бы в себе локальную простоту многочлена невысокой степени и глобальную на всем отрезке $[a, b]$ гладкость, привела к появлению в 1946 г. так называемых *сплайн-функций* или *сплайнов* — специальным образом построенных гладких кусочно-многочленных функций. Получив в 60-х годах XX в. распространение как средство интерполяции сложных кривых, сплайны к настоящему времени стали важной составной частью самых различных вычислительных методов и нашли широчайшее применение в решении разнообразных научно-технических и инженерных задач.

Дадим строгое определение сплайна. Пусть отрезок $[a, b]$ разбит точками, $a = x_0 < x_1 < \dots < x_n = b$ на n частичных отрезков $[x_{i-1}, x_i]$. Сплайном степени m называется функция $S_m(x)$, обладающая следующими свойствами:

- 1) функция $S_m(x)$ непрерывна на отрезке $[a, b]$ вместе со всеми своими производными $S_m^{(1)}(x), S_m^{(2)}(x), \dots, S_m^{(p)}(x)$ до некоторого порядка p ;
- 2) на каждом частичном отрезке $[x_{i-1}, x_i]$ функция $S_m(x)$ совпадает с некоторым алгебраическим многочленом $P_{m,i}(x)$ степени m .

Разность $m - p$ между степенью сплайна и наивысшим порядком непрерывной на отрезке $[a, b]$ производной называется *дефектом сплайна*.

Простейший пример сплайна дает непрерывная кусочно-линейная функция (рис. 11.8), являющаяся сплайном первой степени (*линейным сплайном*) с дефектом, равным 1. Действительно, на отрезке $[a, b]$ сама функция $S_1(x)$ (нулевая производная) непрерывна. В то же время на каждом частичном отрезке $S_1(x)$ совпадает с некоторым многочленом первой степени.

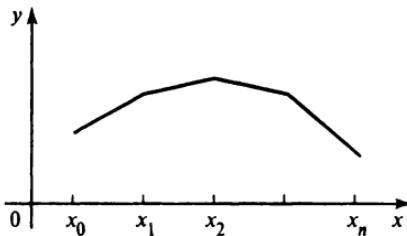


Рис. 11.8

Наиболее широкое распространение на практике получили сплайны $S_3(x)$ третьей степени (*кубические сплайны*) с дефектом, равным 1 или 2. Такие сплайны на каждом из частичных отрезков $[x_{i-1}, x_i]$ совпадают с кубическим многочленом:

$$S_3(x) = P_{3,i}(x) = a_i + b_i(x - x_{i-1}) + c_i(x - x_{i-1})^2 + d_i(x - x_{i-1})^3 \quad (11.63)$$

и имеют на отрезке $[a, b]$ по крайней мере одну непрерывную производную $S'_3(x)$.

Термин «сплайн» происходит от английского слова «spline» (гибкая линейка, стержень) — названия приспособления, использовавшегося чертежниками для проведения гладких кривых через заданные точки. Если гибкую стальную линейку поставить на ребро и, изогнув, зафиксировать ее положение в узловых точках (рис. 11.9), то получится механи-

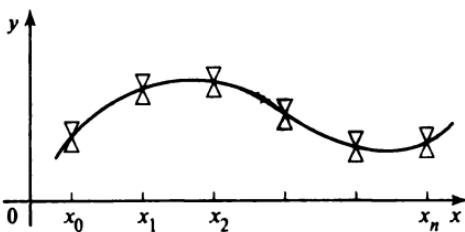


Рис. 11.9

ческий аналог кубического сплайна. Из курса сопротивления материалов известно, что уравнение свободного равновесия профиля $S(x)$ линейки таково: $S^{(4)}(x) = 0$. Следовательно, в промежутке между двумя соседними узлами $S(x)$ представляет собой многочлен третьей степени. В то же время отсутствие у линейки изломов свидетельствует о непрерывности касательной к графику функции $S(x)$ и кривизны, т.е. непрерывности производных $S'(x)$ и $S''(x)$.

2. Интерполяционный сплайн. Пусть функция $y = f(x)$ задана таблицей своих значений $y_i = f(x_i)$, $i = 0, 1, \dots, n$. Сплайн $S_m(x)$ называется **интерполяционным**, если $S_m(x_i) = y_i$ для всех $i = 0, 1, \dots, n$. Значение $s_i = S'_m(x_i)$ называется **наклоном сплайна** в точке x_i .

Заметим, что на отрезке $[x_{i-1}, x_i]$ интерполяционный кубический сплайн однозначно определяется заданием значений y_{i-1} , y_i , s_{i-1} , s_i . В самом деле, из равенства (11.31) вытекает следующая формула:

$$\begin{aligned} S_3(x) = P_{3,i}(x) &= \frac{(x-x_i)^2(2(x-x_{i-1})+h_i)}{h_i^3} y_{i-1} + \frac{(x-x_{i-1})^2(2(x_i-x)+h_i)}{h_i^3} y_i + \\ &+ \frac{(x-x_i)^2(x-x_{i-1})}{h_i^2} s_{i-1} + \frac{(x-x_{i-1})^2(x-x_i)}{h_i^2} s_i, \end{aligned} \quad (11.64)$$

где $h_i = x_i - x_{i-1}$.

Формулу (11.64) можно записать более компактно

$$\begin{aligned} S_3(x_{i-1} + th_i) &= P_{3,i}(x_{i-1} + th_i) = \\ &= (1-t)^2(2t+1)y_{i-1} + t^2(3-2t)y_i + (1-t)^2th_is_{i-1} + t^2(t-1)h_is_i, \end{aligned}$$

используя локальную переменную $t = (x - x_i)/h_i$.

Различные методы интерполяции кубическими сплайнами отличаются один от другого способом выбора наклонов s_i . Обсудим некоторые из них.

3. Локальный сплайн. Если в точках x_i известны значения производной $y'_i = f'(x_i)$, то естественно положить $s_i = y'_i$ для всех $i = 0, 1, \dots, n$. Тогда на каждом частичном отрезке $[x_{i-1}, x_i]$ в соответствии с формулой (11.64) сплайн однозначно определяется значениями y_{i-1} , y_i , y'_{i-1} , y'_i (поэтому его и называют *локальным сплайном*). Заметим, что он совпадает с кубическим интерполяционным многочленом Эрмита (11.31) для отрезка $[x_{i-1}, x_i]$.

Из неравенства (11.33) получается следующая оценка погрешности интерполяции локальным кубическим сплайном:

$$\max_{[a, b]} |f(x) - S_3(x)| \leq \frac{M_4}{384} h_{\max}^4, \quad (11.65)$$

где $h_{\max} = \max_{1 \leq i \leq n} h_i$ — максимальная из длин частичных отрезков.

Заметим, что для построенного указанным образом сплайна можно гарантировать непрерывность на отрезке $[a, b]$ только функции S_3 и ее первой производной S'_3 , т.е. его дефект равен 2.

Существуют и другие способы выбора коэффициентов s_i , приводящие к локальным сплайнам (кубический многочлен Бесселя¹, метод Акимы и др. [18]).

4. Глобальные способы построения кубических сплайнов. Для того чтобы сплайн $S_3(x)$ имел непрерывную на отрезке $[a, b]$ вторую производную $S''_3(x)$, необходимо выбирать наклоны s_i так, чтобы в точках x_i «стыка» многочленов $P_{3, i}$ и $P_{3, i+1}$ совпадали значения их вторых производных:

$$P''_{3, i}(x_i) = P''_{3, i+1}(x_i), \quad i = 1, 2, \dots, n-1. \quad (11.66)$$

Пользуясь формулой (11.64), находим значение

$$P''_{3, i}(x_i) = \frac{2s_{i-1}}{h_i} + \frac{4s_i}{h_i} - 6 \frac{y_i - y_{i-1}}{h_i^2}. \quad (11.67)$$

Из подобной формулы, записанной для многочлена $P_{3, i+1}$, имеем

$$P''_{3, i+1}(x_i) = -\frac{4s_i}{h_{i+1}} - \frac{2s_{i+1}}{h_{i+1}} + 6 \frac{y_{i+1} - y_i}{h_{i+1}^2}. \quad (11.68)$$

¹Фридрих Вильгельм Бессель (1784—1846) — немецкий астроном.

Таким образом, равенства (11.66) приводят к следующей системе уравнений относительно коэффициентов s_i :

$$h_i^{-1}s_{i-1} + 2(h_i^{-1} + h_{i+1}^{-1})s_i + h_{i+1}^{-1}s_{i+1} = 3[h_i^{-2}(y_i - y_{i-1}) + h_{i+1}^{-2}(y_{i+1} - y_i)], \\ i = 1, 2, \dots, n-1. \quad (11.69)$$

Заметим, что эта система уравнений недоопределенна, так как число уравнений системы (равное $n-1$) меньше числа неизвестных (равного $n+1$). Выбор двух оставшихся уравнений обычно связывают с некоторыми дополнительными условиями, накладываемыми на сплайн в граничных точках a и b (*граничными условиями*). Укажем на некоторые из наиболее известных граничных условий.

1°. Если в граничных точках известны значения первой производной $f'(a)$ и $f'(b)$, то естественно положить

$$s_0 = f'(a), \quad s_n = f'(b). \quad (11.70)$$

Дополняя систему (11.69) уравнениями (11.70), приходим к системе уравнений с трехдиагональной матрицей, которая легко решается методом прогонки (см. гл. 5). Полученный таким образом сплайн называется *фундаментальным кубическим сплайном*.

2°. Если в граничных точках известны значения второй производной $f''(a)$ и $f''(b)$, то можно наложить на сплайн граничные условия $S_3''(a) = P_{3,1}''(x_0) = f''(a)$, $S_3''(b) = P_{3,n}''(x_n) = f''(b)$, что приводит к следующим уравнениям:

$$-\frac{4s_0}{h_1} - \frac{2s_1}{h_1} + 6\frac{y_1 - y_0}{h_1^2} = f''(a), \quad (11.71)$$

$$\frac{2s_{n-1}}{h_n} + \frac{4s_n}{h_n} - 6\frac{y_n - y_{n-1}}{h_n^2} = f''(b) \quad (11.72)$$

(достаточно в равенстве (11.68) взять $i=0$, а в равенстве (11.67) $i=n$).

3°. Полагая в уравнениях (11.71), (11.72) $f''(a) = 0, f''(b) = 0$ (независимо от того, выполнены ли эти условия для интерполируемой функции), приходим к системе уравнений, определяющих так называемый *естественный кубический сплайн*.

4°. Часто нет никакой дополнительной информации о значениях производных на концах отрезка. Один из применяемых в этой ситуации подходов состоит в использовании условия «отсутствия узла». Выбор наклонов s_i производят таким образом, чтобы для получаемого сплайна выполнялись условия $P_{3,1}(x) \equiv P_{3,2}(x)$, $P_{3,n-1}(x) \equiv P_{3,n}(x)$. Для этого

достаточно потребовать совпадения в точках x_1 и x_{n-1} соответствующих третьих производных:

$$P_{3,1}^{(3)}(x_1) = P_{3,2}^{(3)}(x_1), \quad P_{3,n-1}^{(3)}(x_{n-1}) = P_{3,n}^{(3)}(x_{n-1}).$$

Эквивалентные алгебраические уравнения выглядят так:

$$2h_1^{-3}(y_0 - y_1) + h_1^{-2}(s_0 + s_1) = 2h_2^{-3}(y_1 - y_2) + h_2^{-2}(s_1 + s_2), \quad (11.73)$$

$$\begin{aligned} 2h_{n-1}^{-3}(y_{n-2} - y_{n-1}) + h_{n-1}^{-2}(s_{n-2} + s_{n-1}) &= \\ &= 2h_n^{-3}(y_{n-1} - y_n) + h_n^{-2}(s_{n-1} + s_n). \end{aligned} \quad (11.74)$$

Та же аппроксимирующая функция может быть получена несколько иначе. Уменьшим число частичных отрезков, объединив попарно отрезки $[x_0, x_1]$, $[x_1, x_2]$ и $[x_{n-2}, x_{n-1}]$, $[x_{n-1}, x_n]$. Это отвечает разбиению отрезка $[a, b]$ точками $a = \tilde{x}_0 < \tilde{x}_1 < \dots < \tilde{x}_{n-2} = b$, где $\tilde{x}_i = x_{i+1}$ для $i = 1, 2, \dots, n-3$, и построению соответствующего интерполяционного сплайна $\tilde{S}_3(x)$. Условия «отсутствия узла» эквивалентны требованию совпадения значений сплайна $\tilde{S}_3(x)$ в точках x_1 и x_{n-1} со значениями y_1 и y_{n-1} .

5°. Если f — периодическая функция с периодом, равным $b - a$, то систему (11.69) следует дополнить условиями

$$s_0 = s_n,$$

$$h_n^{-1}(s_{n-1} + 2s_n) + h_1^{-1}(2s_0 + s_1) = 3[h_n^{-2}(y_n - y_{n-1}) + h_1^{-2}(y_1 - y_0)].$$

Существуют и другие подходы к заданию граничных условий (подробнее об этом см. [18]).

Пример 11.13. Для функции, заданной табл. 11.12, построим (естественный) кубический сплайн. В этом случае система уравнений для наклонов) s_0, s_1, \dots, s_4 в точках x_0, x_1, \dots, x_4 записывается следующим образом:

$$s_0 = -0.5s_1 + 1.2, \quad (11.75)$$

$$s_0 + 4s_1 + s_2 = 3.6, \quad (11.76)$$

$$s_1 + 4s_2 + s_3 = -1.2, \quad (11.77)$$

$$s_2 + 4s_3 + s_4 = -3.6, \quad (11.78)$$

$$s_4 = -0.5s_3 - 0.6. \quad (11.79)$$

Решая ее, получаем значения $s_0 = 57/70$, $s_1 = 54/70$, $s_2 = -3/10$, $s_3 = -54/70$, $s_4 = -3/14$. Теперь на каждом частичном отрезке значения

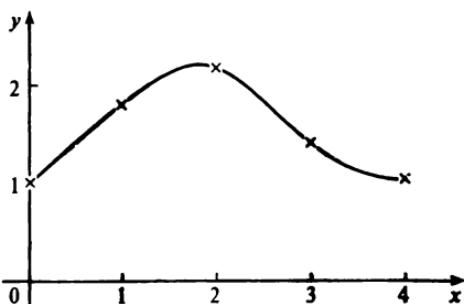


Рис. 11.10

сплайна можно вычислить по формуле (11.64). Соответствующий график приведен на рис. 11.10 (ср. с рис. 11.6 и 11.7). ▲

Пример 11.14. Интерполируем функцию, заданную табл. 11.12, кубическим сплайном, используя условие «отсутствия узла». В этом случае уравнения (11.76)–(11.78) останутся прежними, а уравнения (11.75) и (11.79) заменяются следующими:

$$s_0 - s_2 = 0.8, \quad s_2 - s_4 = -0.8. \quad (11.80)$$

Решая систему (11.76)–(11.78), (11.80), получаем значения $s_0 = 8/15$, $s_1 = 5/6$, $s_2 = -4/15$, $s_3 = -29/30$, $s_4 = 8/15$. График соответствующего сплайна мало отличается от графика, изображенного на рис. 11.10. ▲

5. Погрешность приближения кубическими сплайнами

Теорема 11.9. Пусть функция f имеет на отрезке $[a, b]$ непрерывную производную четвертого порядка и $M_4 = \max_{[a, b]} |f^{(4)}(x)|$. Тогда для интерполяционного кубического сплайна $S_3(x)$, удовлетворяющего граничным условиям типов 1° , 2° , 4° или 5° (последнее — для случая периодической функции), справедлива следующая оценка погрешности:

$$\max_{[a, b]} |f(x) - S_3(x)| \leq CM_4 h_{\max}^4. \blacksquare \quad (11.81)$$

Заметим, что сплайн S_3 не только сам аппроксимирует функцию f , но его производные S'_3 , S''_3 , $S^{(3)}_3$ приближают соответствующие производные функции f . Сформулируем соответствующую теорему в наиболее простом случае, когда таблица задана с постоянным шагом h .

Теорема 11.10. При выполнении условий теоремы 11.9 для указанных в ней сплайнов справедливы неравенства

$$\max_{[a, b]} |f^{(k)}(x) - S_3^{(k)}(x)| \leq C_k M_4 h^{4-k}, \quad k = 1, 2, 3. \blacksquare \quad (11.82)$$

Замечание 1. Естественный кубический сплайн обладает следующим замечательным минимальным свойством. Он дает минимум функционалу «энергии»

$$J(g) = \int_a^b |g''(x)|^2 dx$$

на множестве всех функций g , которые непрерывно дифференцируемы на отрезке $[a, b]$, имеют кусочно-непрерывную на отрезке $[a, b]$ вторую производную и удовлетворяют условиям $g(x_i) = y_i$, $0 \leq i \leq n$.

Замечание 2. Благодаря большей простоте записи, указанному выше минимальному свойству и благозвучному названию естественные сплайны получили значительное распространение. Однако искусственное наложение условий $f''(a) = 0$, $f''(b) = 0$ при интерполяции функций, которые этим условиям не удовлетворяют, приводит к значительной потере точности. Вместо четвертого порядка точности (как локальный кубический сплайн или кубические сплайны с граничными условиями типов 1° , 2° , 4° , 5°) естественный сплайн обладает лишь вторым порядком точности. Если использование естественного сплайна не вызвано какими-то специальными причинами, то следует, по-видимому, отказаться от него в пользу кубического сплайна с граничным условием типа 4° .

§ 11.12. Понятие о дискретном преобразовании Фурье и тригонометрической интерполяции

1. Дискретное преобразование Фурье. В прикладных исследованиях широко используются различные варианты преобразования Фурье¹ функций непрерывного аргумента, а также представление функций в виде сходящихся тригонометрических рядов (рядов Фурье). Известно, например, что всякая непрерывно дифференцируемая периодическая с периодом 1 функция f может быть разложена в ряд Фурье:

$$f(x) = \sum_{k=-\infty}^{\infty} \alpha_k \exp(2\pi i kx). \quad (11.83)$$

Здесь i — мнимая единица. Коэффициенты разложения вычисляются по формулам

$$\alpha_k = \int_0^1 f(x) \exp\{-2\pi i kx\} dx. \quad (11.84)$$

¹ Жан Батист Жозеф Фурье (1768—1830) — французский математик, один из основоположников математической физики

Однако нередко функция f бывает задана лишь в конечном числе точек $x_j = j/N$, $j = 0, 1, \dots, N - 1$. В этом случае аналогом формулы (11.83) является разложение вида

$$f(x_j) = \sum_{k=0}^{N-1} a_k \exp\{2\pi i k x_j\}, \quad 0 \leq j < N. \quad (11.85)$$

Заметим, что это разложение имеет место тогда и только тогда, когда тригонометрический многочлен

$$S_N(x) = \sum_{k=0}^{N-1} a_k \exp\{2\pi i k x\} \quad (11.86)$$

интерполирует функцию f по ее значениям в точках x_j , $0 \leq j < N$. Выше (см. пример 11.2) было доказано, что система функций $\phi_k(x) = \exp\{2\pi i k x\}$, $0 \leq k < N$ ортогональна на множестве точек $x_j = j/N$, $0 \leq j < N$, причем $(\phi_k, \phi_k) = N$. Следовательно, разложение (11.85) действительно имеет место, причем в силу равенства (11.19) коэффициенты a_k определяются по формуле

$$a_k = \frac{1}{N} \sum_{l=0}^{N-1} f(x_l) \exp\{-2\pi i k x_l\}, \quad 0 \leq k < N. \quad (11.87)$$

Операцию преобразования набора значений $f(x_0), f(x_1), \dots, f(x_{N-1})$ в набор коэффициентов a_0, a_1, \dots, a_{N-1} принято называть *прямым дискретным преобразованием Фурье*, а обратную операцию — *обратным дискретным преобразованием Фурье*. Осуществление этих операций является важной составной частью многих алгоритмов.

Для удобства изложения введем обозначение $\omega = \exp\{2\pi i / N\}$ и перепишем формулы (11.85), (11.87) в следующем виде:

$$f(x_j) = \sum_{k=0}^{N-1} a_k \omega^{kj}, \quad 0 \leq j < N, \quad (11.88)$$

$$a_k = \frac{1}{N} \sum_{l=0}^{N-1} f(x_l) \omega^{-kl}, \quad 0 \leq k < N. \quad (11.89)$$

2. Быстрое дискретное преобразование Фурье. Если вычисления проводить непосредственно по формулам (11.88) и (11.89), то на выполнение каждого из преобразований потребуется примерно N^2 арифметических операций. (Здесь под арифметической операцией понимается умножение двух комплексных чисел с последующим сложением. Величины ω^{kj} считаются вычисленными заранее.)

Однако, когда число N не является простым, количество арифметических операций, требуемых для вычисления по формулам (11.88) и (11.89), можно существенно уменьшить. Поясним сказанное на примере вычислений по формулам (11.88). (Вычисления по формулам (11.89) производятся аналогично с заменой ω на ω^{-1} .)

Пусть $N = N_1 N_2$, где $2 \leq N_1, 2 \leq N_2$ — целые числа. Представим индекс j виде $j = j_1 N_1 + j_0$, где $0 \leq j_1 < N_2, 0 \leq j_0 < N_1$. Положим $k = k_1 N_2 + k_0$, где $0 \leq k_1 < N_1, 0 \leq k_0 < N_2$. Пользуясь тем, что $kj = k_1 j_1 N_1 + k_1 j_0 N_2 + k_0 j$ и $\omega^{k_1 j_1 N_1} = 1$, имеем $\omega^{kj} = \omega^{k_1 j_0 N_2} \omega^{k_0 j}$. Заменяя в формуле (11.88) суммирование по индексу k операцией повторного суммирования по индексам k_0 и k_1 , получаем

$$f(x_j) = \sum_{k_0=0}^{N_2-1} \sum_{k_1=0}^{N_1-1} a_{k_1 N_2 + k_0} \omega^{k_1 j_0 N_2} \omega^{k_0 j} = \sum_{k_0=0}^{N_2-1} \tilde{a}(k_0, j_0) \omega^{k_0 j}, \quad (11.90)$$

где

$$\tilde{a}(k_0, j_0) = \sum_{k_1=0}^{N_1-1} a_{k_1 N_2 + k_0} \omega^{k_1 j_0 N_2}. \quad (11.91)$$

Массив \tilde{a} содержит N чисел и для его вычисления требуется NN_1 арифметических операций. После того как найдены значения $\tilde{a}(k_0, j_0)$, на вычисления по формуле (11.90) требуется NN_2 операций. Таким образом, общее число арифметических операций равно $N(N_1 + N_2)$. Заметим, что достигнута экономия в числе операций, поскольку $N_1 + N_2 < N_1 N_2 = N$, как только $N > 4$. Выигрыша удалось достичь благодаря тому, что оказалось возможным выделить группы слагаемых (11.91), которые используются для вычисления значений $f(x_j)$ при различных j , но сами вычисляются лишь однажды.

Указанная выше идея развита в алгоритмах *быстрого дискретного преобразования Фурье* (БПФ)¹. Если $N = N_1 \cdot N_2 \cdots N_m$ (где $2 \leq N_s$), то с помощью БПФ можно выполнить дискретное преобразование Фурье за $N(N_1 + N_2 + \dots + N_m)$ арифметических операций. Особенно эффективным является этот алгоритм, когда число N является степенью числа 2

¹ Наиболее раннее описание БПФ появилось в 1866 г. и принадлежит Гауссу. В течение следующих 100 лет алгоритм был переоткрыт несколько раз. Отправной точкой современного применения БПФ послужила вышедшая в 1965 г. работа Джона Кули и Джона Тьюки, поэтому иногда этот алгоритм называют *алгоритмом Кули—Тьюки*.

Интересную историю открытия БПФ и полезное обсуждение различных его аспектов можно найти в [54], см. также [9].

(т.е. $N = 2^m$). В этом случае вместо N^2 операций требуется выполнить лишь $2N \log_2 N$ операций. Например, для $N = 1024 = 2^{10}$ этот алгоритм позволяет ускорить вычисления в $N/(2 \log_2 N) = 1024/20 \approx 50$ раз.

Замечание 1. В случае, когда число N не является степенью двойки, также существует возможность выполнения БПФ за $O(N \log_2 N)$ арифметических операций [12*].

Широкое внедрение алгоритма БПФ в практические вычисления привело к подлинной революции во многих областях, связанных с обработкой числовой информации. Программы, реализующие различные варианты этого алгоритма, входят в стандартное математическое обеспечение компьютеров и доступны массовому пользователю.

Замечание 2. Часто разложение (11.85) записывают в эквивалентном виде:

$$f(x_j) = \sum_{-N/2 < k \leq N/2} a_k \exp\{2\pi i k x_j\},$$

что соответствует интерполяции тригонометрическим многочленом

$$S_N(x) = \sum_{-N/2 < k \leq N/2} a_k \exp\{2\pi i k x\}. \quad (11.92)$$

Здесь коэффициенты a_k по-прежнему задаются формулой (11.87), но для $-N/2 < k \leq N/2$.

Замечание 3. Хотя интерполяционные тригонометрические многочлены (11.86), (11.92) и совпадают в точках x_j , они принимают существенно разные значения в точках x , отличных от узловых. Для интерполяции предпочтительней формула (11.92).

3. Тригонометрическая интерполяция. Рассмотрим кратко задачу интерполяции функции f , заданной в точках $0 \leq x_0 < x_1 < \dots < x_{N-1} \leq 1$ тригонометрическим многочленом (11.92). К ней приводят, например типичная радиотехническая задача о тригонометрической интерполяции периодического сигнала.

Не вдаваясь в довольно сложную проблему оценки погрешности тригонометрической интерполяции, отметим тем не менее, что для гладкой периодической с периодом 1 функции f есть основание рассчитывать на выполнение приближенного равенства $f(x) \approx S_N(x)$ для всех $x \in [0, 1]$.

Рассмотрим важный вопрос о чувствительности многочлена S_N к погрешностям в исходных данных. Пусть значения $y_i^* \approx f(x_i)$ интерполируемой функции задаются с погрешностями ε_i и известно, что

$|\varepsilon_i| \leq \bar{\Delta}(y^*)$ для $i = 0, 1, \dots, N - 1$. Тогда вычисляемый по значениям y_i^* тригонометрический интерполяционный многочлен S_N^* содержит погрешность. Для нее справедлива оценка

$$\bar{\Delta}(S_N^*) = \max_{[0, 1]} |S_N^*(x) - S_N(x)| \leq \tilde{\Lambda}_N \bar{\Delta}(y^*),$$

аналогичная оценке (11.61) для алгебраических многочленов. Здесь $\tilde{\Lambda}_N$ — постоянная, являющаяся аналогом константы Лебега Λ_N .

Примечательно то, что в отличие от задачи интерполяции алгебраическими многочленами (см. § 11.10) оптимальным (т.е. дающим минимальное значение $\tilde{\Lambda}_N$) является равномерное распределение узлов, которому отвечает значение $\tilde{\Lambda}_N \approx \frac{2}{\pi} \ln[(N+1)/2]$.

Таким образом, при тригонометрической интерполяции выбор узлов $x_j = j/N$ ($0 \leq j < N$) является наиболее естественным с точки зрения как простоты вычисления коэффициентов многочлена (быстрое дискретное преобразование Фурье), так и минимизации влияния ошибок исходных данных.

§ 11.13. Метод наименьших квадратов

Задача наименьших квадратов возникает в самых различных областях науки и техники. Например, к ней приходят при статистической обработке экспериментальных данных с помощью регрессионного анализа. В инженерной деятельности задача наименьших квадратов используется в таких областях, как оценивание параметров и фильтрация.

1. Линейная задача наименьших квадратов. Пусть функция $y = f(x)$ задана таблицей приближенных значений

$$y_i \approx f(x_i), \quad i = 0, 1, \dots, n, \quad (11.93)$$

полученных с ошибками $\varepsilon_i = y_i^0 - y_i$, где $y_i^0 = f(x_i)$. Если значения y_i получены из эксперимента, то ошибки носят случайный характер и часто уровень погрешности («шума» таблицы) бывает значительным (рис. 11.11, а).

Предположим, что для аппроксимации функции f используется линейная модель:

$$y = \Phi_m(x) \equiv a_0 \varphi_0(x) + a_1 \varphi_1(x) + \dots + a_m \varphi_m(x), \quad (11.94)$$

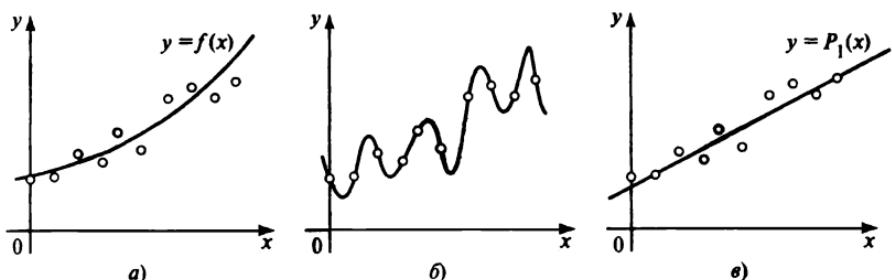


Рис. 11.11

здесь $\phi_0(x), \phi_1(x), \dots, \phi_m(x)$ — заданные базисные функции¹; a_0, a_1, \dots, a_m — параметры модели, являющиеся коэффициентами обобщенного многочлена $\Phi_m(x)$.

Как уже отмечалось выше, одной из наиболее простых и часто используемых линейных моделей вида (11.94) (при $\phi_k(x) = x^k$) является полиномиальная модель

$$y = P_m(x) \equiv a_0 + a_1x + \dots + a_mx^m. \quad (11.95)$$

Когда уровень неопределенности исходных данных высок, неестественно требовать от модели (11.94) выполнения условий (11.7) совпадения значений обобщенного многочлена $\Phi_m(x)$ в точках x_i , с заданными значениями y_i , т.е. использовать интерполяцию. Как нетрудно видеть (см. рис. 11.11, б), при интерполировании происходит повторение ошибок наблюдений, в то время как при обработке экспериментальных данных желательно, напротив, их сглаживание.

Отказываясь от требования выполнения в точках x_i точных равенств (11.7), следует все же стремиться к тому, чтобы в этих точках выполнялись соответствующие приближенные равенства

$$\begin{aligned} a_0\phi_0(x_0) + a_1\phi_1(x_0) + \dots + a_m\phi_m(x_0) &\approx y_0, \\ a_0\phi_0(x_1) + a_1\phi_1(x_1) + \dots + a_m\phi_m(x_1) &\approx y_1, \\ \dots \\ a_0\phi_0(x_n) + a_1\phi_1(x_n) + \dots + a_m\phi_m(x_n) &\approx y_n. \end{aligned} \quad (11.96)$$

¹ В данном параграфе рассматриваются функции, принимающие только вещественные значения.

Используя обозначения (11.9), запишем систему приближенных равенств (11.96) в матричном виде:

$$\mathbf{P}\mathbf{a} \approx \mathbf{y}. \quad (11.97)$$

Из различных критериев, позволяющих выбрать параметры a_0, a_1, \dots, a_m , модели (11.94) так, чтобы приближенные равенства (11.96) удовлетворялись наилучшим в некотором смысле образом, наиболее часто используется критерий наименьших квадратов. Согласно этому критерию параметры выбираются так, чтобы минимизировать *среднеквадратичное уклонение*

$$\delta(\Phi_m, \mathbf{y}) = \sqrt{\frac{1}{n+1} \sum_{i=0}^n (\Phi_m(x_i) - y_i)^2}$$

обобщенного многочлена $\Phi_m(x) = \sum_{j=0}^m a_j \varphi_j(x)$ от заданных табличных значений y_i ($0 \leq i \leq n$). Заметим, что минимум среднеквадратичного уклонения достигается при тех же значениях a_0, a_1, \dots, a_m , что и минимум функции

$$s(\mathbf{a}, \mathbf{y}) = \sum_{i=0}^n \left(\sum_{j=0}^m a_j \varphi_j(x_i) - y_i \right)^2 = \|\mathbf{P}\mathbf{a} - \mathbf{y}\|_2^2,$$

причем

$$\delta(\Phi_m, \mathbf{y})^2 = \frac{1}{n+1} s(\mathbf{a}, \mathbf{y}).$$

Итак, линейная задача метода наименьших квадратов состоит в следующем. Требуется найти (при фиксированном наборе функций $\varphi_0, \varphi_1, \dots, \varphi_m$) обобщенный многочлен $\Phi_m^*(x)$, для которого среднеквадратичное уклонение принимает минимальное значение:

$$\delta(\Phi_m^*, \mathbf{y}) = \min_{\Phi_m} \delta(\Phi_m, \mathbf{y}).$$

Искомый обобщенный многочлен Φ_m^* будем далее называть *многочленом наилучшего среднеквадратичного приближения*.

2. Нормальная система. Существуют различные подходы к решению поставленной задачи. Простейший из них состоит в использовании необходимого условия экстремума функции s :

$$\frac{\partial s}{\partial a_k} = 0, \quad k = 0, 1, \dots, m. \quad (11.98)$$

Вычисляя частные производные функции s и изменяя порядок суммирования, от равенств (11.98) переходим к системе линейных алгебраических уравнений

$$\sum_{j=0}^m \left(\sum_{i=0}^n \varphi_j(x_i) \varphi_k(x_i) \right) a_j = \sum_{i=0}^n y_i \varphi_k(x_i), \quad k = 0, 1, \dots, m. \quad (11.99)$$

которая называется *нормальной системой метода наименьших квадратов*.

Как нетрудно видеть, нормальную систему можно записать в виде

$$P^T P a = P^T y, \quad (11.100)$$

или, используя матрицу Грама $\Gamma = P^T P$ (см. § 11.2) и вводя вектор $b = P^T y$, в виде

$$\Gamma a = b. \quad (11.101)$$

Лемма 11.1. Пусть a — решение системы (11.101). Тогда для любого $a' = a + \Delta a$ имеет место равенство

$$s(a', y) = s(a, y) + \|P\Delta a\|_2^2.$$

Доказательство. Имеем

$$\begin{aligned} s(a', y) &= \|P(a + \Delta a) - y\|_2^2 = \|Pa - y\|_2^2 + 2(Pa - y, P\Delta a) + \|P\Delta a\|_2^2 = \\ &= s(a, y) + 2(P^T Pa - P^T y, \Delta a) + \|P\Delta a\|_2^2. \end{aligned}$$

Остается заметить, что второе слагаемое в правой части полученного равенства в силу (11.100) равно нулю. ■

Теорема 11.11. Пусть система функций $\varphi_0, \varphi_1, \dots, \varphi_m$ линейно независима в точках x_0, x_1, \dots, x_n . Тогда многочлен наилучшего среднеквадратичного приближения Φ_m^y существует и единственен.

Доказательство. В силу теоремы 11.1 $\det \Gamma \neq 0$. Поэтому решение a системы (11.101) существует и единственno. Таким образом, если многочлен Φ_m^y существует, то его коэффициенты a_0, a_1, \dots, a_m определяются единственным образом.

Заметим теперь, что согласно лемме 11.1 для любого $a' \neq a$ имеем $s(a', y) \geq s(a, y)$, т.е. на решении a нормальной системы действительно достигается минимум функции s . ■

Замечание 1. Если $m = n$ и система функций $\Phi_0, \Phi_1, \dots, \Phi_n$ линейно независима в точках x_0, x_1, \dots, x_n , то многочлен Φ_n^y , найденный методом наименьших квадратов, совпадает с интерполяционным многочленом Φ_n . В самом деле, $\Phi_n(x_i) = y_i$ для всех $i = 0, 1, \dots, n$ и поэтому $\delta(\Phi_n, y) = 0$. Так как среднеквадратичное уклонение не может быть отрицательным, то Φ_n — многочлен наилучшего среднеквадратичного приближения. В силу его единственности $\Phi_n = \Phi_n^y$.

Замечание 2. Как правило, при использовании метода наименьших квадратов предполагается что $m \ll n$. В этом случае метод обладает некоторыми сглаживающими свойствами.

Очень часто для приближения по методу наименьших квадратов используются алгебраические многочлены степени $m \leq n$. Поскольку система функций $1, x, \dots, x^m$ линейно независима в точках x_0, x_1, \dots, x_n при $m \leq n$ (см. пример 11.1), в силу теоремы 11.11 алгебраический многочлен наилучшего среднеквадратичного приближения существует и единственен.

Так как в случае приближения алгебраическими многочленами $\varphi_k(x) = x^k$, то нормальная система (11.99) принимает следующий вид:

$$\sum_{j=0}^m \left(\sum_{i=0}^n x_i^{j+k} \right) a_j = \sum_{i=0}^n y_i x_i^k, \quad k = 0, 1, \dots, m. \quad (11.102)$$

Запишем систему (11.102) в развернутом виде в двух наиболее простых случаях $m = 1$ и $m = 2$. Когда приближение осуществляется многочленом первой степени $P_1(x) = a_0 + a_1 x$, нормальная система имеет вид:

$$\begin{aligned} (n+1)a_0 + \left(\sum_{i=0}^n x_i \right) a_1 &= \sum_{i=0}^n y_i, \\ \left(\sum_{i=0}^n x_i \right) a_0 + \left(\sum_{i=0}^n x_i^2 \right) a_1 &= \sum_{i=0}^n y_i x_i. \end{aligned} \quad (11.103)$$

Если же используется многочлен второй степени $P_2(x) = a_0 + a_1 x + a_2 x^2$, то нормальная система имеет вид:

$$(n+1)a_0 + \left(\sum_{i=0}^n x_i \right) a_1 + \left(\sum_{i=0}^n x_i^2 \right) a_2 = \sum_{i=0}^n y_i,$$

$$\left(\sum_{i=0}^n x_i \right) a_0 + \left(\sum_{i=0}^n x_i^2 \right) a_1 + \left(\sum_{i=0}^n x_i^3 \right) a_2 = \sum_{i=0}^n y_i x_i, \quad (11.104)$$

$$\left(\sum_{i=0}^n x_i^2 \right) a_0 + \left(\sum_{i=0}^n x_i^3 \right) a_1 + \left(\sum_{i=0}^n x_i^4 \right) a_2 = \sum_{i=0}^n y_i x_i^2.$$

Пример 11.15. Пусть функция $y = f(x)$ задана следующей таблицей (табл. 11.13):

Таблица 11.13

x	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
y	0.21	0.23	0.31	0.29	0.42	0.35	0.58	0.61	0.59	0.66

Используя метод наименьших квадратов, аппроксимируем ее многочленами первой и второй степени и найдем соответствующие среднеквадратичные уклонения δ_1 и δ_2 .

Вычислим коэффициенты и правые части нормальных систем (11.103), (11.104):

$$\sum_{i=0}^9 x_i = 4.5, \quad \sum_{i=0}^9 x_i^2 = 2.85, \quad \sum_{i=0}^9 x_i^3 = 2.025, \quad \sum_{i=0}^9 x_i^4 = 1.5333,$$

$$\sum_{i=0}^9 y_i = 4.25, \quad \sum_{i=0}^9 y_i x_i = 2.356, \quad \sum_{i=0}^9 y_i x_i^2 = 1.6154.$$

Для многочлена первой степени нормальная система имеет вид

$$10a_0 + 4.5a_1 = 4.25,$$

$$4.5a_0 + 2.85a_1 = 2.356.$$

Решив ее, получим значения $a_0 \approx 0.183$, $a_1 \approx 0.538$ коэффициентов многочлена $P_1(x) = a_0 + a_1 x$ наилучшего среднеквадратичного приближения. Его график изображен на рис. 11.11, в.

Запишем теперь нормальную систему для многочлена второй степени:

$$10a_0 + 4.5a_1 + 2.85a_2 = 4.25,$$

$$4.5a_0 + 2.85a_1 + 2.025a_2 = 2.356,$$

$$2.85a_0 + 2.025a_1 + 1.5333a_2 = 1.6154.$$

Решив ее, получим значения $a_0 \approx 1.194$, $a_1 \approx 0.452$, $a_2 \approx 0.0947$ коэффициентов многочлена $P_2(x) = a_0 + a_1x + a_2x^2$ наилучшего среднеквадратичного приближения.

Вычисления по формуле

$$\delta_m = \sqrt{\frac{1}{10} \sum_{i=0}^9 (P_m(x_i) - y_i)^2}$$

для $m = 1$ и $m = 2$ дают значения $\delta_1 \approx 0.0486$, $\delta_2 \approx 0.0481$.

Так как средняя погрешность ϵ исходных данных заведомо превышает 0.01, нетрудно заключить, что приближения многочленами первой и второй степени дают в данной ситуации практически эквивалентный результат. Учитывая большую простоту использования линейных функций, достаточно, по-видимому, остановиться на приближении $f(x) \approx 0.183 + 0.538x$. ▲

3. Некоторые вычислительные аспекты задачи наименьших квадратов. Метод вычисления параметров a_0, a_1, \dots, a_m с помощью решения нормальной системы (11.101) кажется весьма привлекательным. Действительно, задача сводится к стандартной проблеме линейной алгебры — решению системы линейных алгебраических уравнений с квадратной матрицей. Более того, можно показать, что, когда функции $\Phi_0, \Phi_1, \dots, \Phi_m$ линейно независимы в точках x_0, x_1, \dots, x_n , матрица системы Γ является симметричной и положительно определенной. В частности, это означает, что при решении нормальной системы методом Гаусса не нужен выбор главных элементов; возможно также использование метода Холецкого (см. гл. 5).

Следует тем не менее обратить серьезное внимание на то обстоятельство, что при отсутствии специального выбора базисных функций $\Phi_0, \Phi_1, \dots, \Phi_m$ уже при $m \geq 5$ нормальная система обычно оказывается очень плохо обусловленной. Казалось бы, из теоремы 11.11 следует, что единственное ограничение на систему базисных функций состоит в том, что она должна быть линейно независима в заданных точках. Однако, будучи формально линейно независимой, система функций $\Phi_0, \Phi_1, \dots, \Phi_m$ может оказаться очень близкой к линейно зависимой. Использование такой «почти линейно зависимой» системы базисных функций делает задачу метода наименьших квадратов плохо обусловленной. При переходе от задачи наименьших квадратов к задаче решения нормальной системы $P^T P \mathbf{a} = P^T \mathbf{y}$ происходит как бы симметризация системы $P \mathbf{a} \approx \mathbf{y}$ (см. § 6.4). При этом еще более ухудшается обусловленность задачи. В этом случае вычисленные на компьютере

как решение системы (11.101) параметры модели могут оказаться полностью искаженными ошибками округления.

Простейший пример такой «почти линейно зависимой» системы базисных функций при больших m дает система $1, x, \dots, x^m$, широко применяемая при аппроксимации алгебраическими многочленами. При $m \geq 5$ соответствующая нормальная система, как правило, здесь настолько плохо обусловлена, что ее использование практически бесполезно.

В определенном смысле «наиболее линейно независимой» является система функций $\varphi_0, \varphi_1, \dots, \varphi_m$, ортогональных на множестве точек x_0, x_1, \dots, x_n . Матрица Грама такой системы диагональна, а потому решение нормальной системы (11.101) вычисляется легко:

$$a_k = b_k / \gamma_{kk}, \quad b_k = \sum_{i=0}^n y_i \varphi_k(x_i), \quad \gamma_{kk} = \sum_{i=0}^n \varphi_k(x_i)^2, \quad k = 0, 1, \dots, m.$$

Хотя выбор ортогональной на множестве точек x_0, x_1, \dots, x_n системы функций и желателен, он далеко не всегда возможен и удобен. Поэтому часто используются системы базисных функций, для которых матрица Грама лишь близка к диагональной. При аппроксимации на отрезке $[-1, 1]$ алгебраическими многочленами степени m пример такой системы дает система многочленов Чебышева $T_0(x), T_1(x), \dots, T_m(x)$. Заметим, кстати, что найденный методом наименьших квадратов многочлен $\tilde{P}_m(x) = \tilde{a}_0 T_0(x) + \tilde{a}_1 T_1(x) + \dots + \tilde{a}_m T_m(x)$ дает лишь иное, отличное от стандартного (11.95) представление многочлена наилучшего среднеквадратичного приближения. Однако задача определения параметров $\tilde{a}_0, \tilde{a}_1, \dots, \tilde{a}_m$ обладает существенно лучшей обусловленностью и поэтому предпочтительнее с вычислительной точки зрения.

Существуют методы решения задачи наименьших квадратов, предваряющие решение нормальной системы численной ортогонализацией системы базисных функций (см., например, [64, 73]). Однако в настоящее время в серьезной вычислительной практике нормальная система, как правило, не используется. Применяются другие, более надежные методы (учитывающие, например, информацию об уровне погрешности данных и относительной точности используемого компьютера). С одним из таких методов, основанном на сингулярном разложении матрицы P , можно познакомиться в [54, 106]; см. также § 5.12 данной книги.

4. Понятие о статистических свойствах метода наименьших квадратов. Пусть значения y , функции f в точках x , определяются в результате эксперимента. Предположим, что ошибки наблюдения

$\varepsilon_i = y_i^0 - y_i$ являются независимыми случайными величинами с нулевым средним значением и дисперсией равной σ^2 , т.е.

$$M[\varepsilon_i] = 0, \quad M[\varepsilon_i^2] = \sigma^2 \quad (i = 0, 1, \dots, n); \quad (11.105)$$

$$M[\varepsilon_i \varepsilon_j] = 0 \quad (i \neq j, \quad i = 0, 1, \dots, n, \quad j = 0, 1, \dots, n). \quad (11.106)$$

Рассмотрим сначала простейший случай, когда измеряемая величина постоянна, т.е. $f(x) = a$. В этой ситуации естественно искать приближение в виде постоянной a_0 . В изложенной выше схеме это соответствует выбору $m = 0$, $\phi_0(x) = 1$ и $\Phi_0(x) = a_0$, а нормальная система превращается в одно линейное уравнение $(n + 1)a_0 = \sum_{i=0}^n y_i$.

Таким образом, оценкой постоянной величины a по методу наименьших квадратов является среднее арифметическое значение

$$a_0 = \frac{1}{n+1} \sum_{i=0}^n y_i, \quad (11.107)$$

т.е. измеренные значения осредняются так, как это принято в статистике.

Пусть $\Delta a_0 = a - a_0$ — ошибка оценки (11.107). Заметим, что

$$\Delta a_0 = \frac{1}{n+1} \sum_{i=0}^n \varepsilon_i, \text{ а потому}$$

$$M[\Delta a_0] = \frac{1}{n+1} \sum_{i=0}^n M[\varepsilon_i] = 0;$$

$$D[\Delta a_0] = \frac{1}{(n+1)^2} \sum_{i=0}^n \sum_{j=0}^n M[\varepsilon_i \varepsilon_j] = \frac{\sigma^2}{n+1}. \quad (11.108)$$

Итак, математическое ожидание ошибки Δa_0 равно нулю, причем, как видно из равенства (11.108), ее дисперсия стремится к нулю при $n \rightarrow \infty$. Аналогично осредняются случайные ошибки с ростом числа наблюдений и в общем случае.

Пусть $\Phi_m^y(x) = \sum_{k=0}^m a_k^0 \phi_k(x)$, $\Phi_m^y(x) = \sum_{k=0}^m a_k \phi_k(x)$ — многочлены наилучшего среднеквадратичного приближения, первый из которых отвечает вектору $y^0 = (y_0^0, y_1^0, \dots, y_n^0)^T$ данных, не содержащих ошибок, а второй вектору $y = (y_0, y_1, \dots, y_n)^T$ данных, содержащих случайные ошибки.

Положим $\mathbf{a}^0 = (a_0^0, a_1^0, \dots, a_m^0)^T$, $\Delta \mathbf{a} = (\Delta a_0, \Delta a_1, \dots, \Delta a_n)^T$ (где $\Delta a_k = a_k^0 - a_k$, $k = 0, 1, \dots, m$) и $\varepsilon = (\varepsilon_0, \varepsilon_1, \dots, \varepsilon_n)^T$. Рассмотрим многочлен

$$\Delta \Phi_m(x) = \Phi_m^{y^0}(x) - \Phi_m^y(x) = \sum_{k=0}^m \Delta a_k \varphi_k(x).$$

Так как векторы \mathbf{a} и \mathbf{a}^0 удовлетворяют системам уравнений $\mathbf{\Gamma a} = \mathbf{P}^T y$, $\mathbf{\Gamma a}^0 = \mathbf{P}^T y^0$, то вектор $\Delta \mathbf{a}$ является решением системы

$$\mathbf{\Gamma} \Delta \mathbf{a} = \mathbf{P}^T \varepsilon. \quad (11.109)$$

Откуда вытекает равенство

$$\Delta \mathbf{a} = \mathbf{\Gamma}^{-1} \mathbf{P}^T \varepsilon. \quad (11.110)$$

Заметим, что $\Delta \Phi_m(x)$ и Δa_k — это случайные ошибки значения многочлена $\Phi_m^y(x)$ и его коэффициентов a_k , вызванные наличием случайных ошибок ε , в исходных данных. Покажем, что эти ошибки имеют нулевое математическое ожидание.

Действительно, из формулы (11.110) следует, что

$$M[\Delta \mathbf{a}] = \mathbf{\Gamma}^{-1} \mathbf{P}^T M[\varepsilon] = 0,$$

откуда в свою очередь имеем

$$M[\Delta \Phi_m(x)] = \sum_{k=0}^m M[\Delta a_k] \varphi_k(x) = 0.$$

Введем величину

$$\rho = \sqrt{\frac{1}{n+1} \sum_{i=0}^n (\Delta \Phi_m(x_i))^2},$$

равную среднеквадратичному значению ошибки $\Delta \Phi_m(x)$.

Теорема 11.12. Справедливо равенство

$$M[\rho^2] = \frac{n+1}{n+1} \sigma^2. \quad (11.111)$$

Доказательство. Заметим, что

$$\begin{aligned} \rho^2 &= \frac{1}{n+1} (\mathbf{P} \Delta \mathbf{a}, \mathbf{P} \Delta \mathbf{a}) = \frac{1}{n+1} (\Delta \mathbf{a}, \mathbf{P}^T \mathbf{P} \Delta \mathbf{a}) = \\ &= \frac{1}{n+1} (\mathbf{\Gamma}^{-1} \mathbf{P}^T \varepsilon, \mathbf{P}^T \varepsilon) = \frac{1}{n+1} (\mathbf{P} \mathbf{\Gamma}^{-1} \mathbf{P}^T \varepsilon, \varepsilon). \end{aligned}$$

Из определения операции умножения матриц следует, что элемент a_{ij} матрицы $A = P\Gamma^{-1}P^T$ имеет вид $a_{ij} = \sum_{l=0}^m \sum_{k=0}^m p_{il} \gamma_{lk}^{(-1)} p_{kj}^T$, где p_{il} , $\gamma_{lk}^{(-1)}$ и p_{kj}^T — элементы матриц P , Γ^{-1} и P^T соответственно. Поэтому

$$\begin{aligned} M[\rho^2] &= \frac{1}{n+1} \sum_{i=0}^n \sum_{j=0}^n a_{ij} M[\varepsilon_i \varepsilon_j] = \frac{1}{n+1} \sum_{i=0}^n a_{ii} \sigma^2 = \\ &= \frac{1}{n+1} \sum_{l=0}^m \sum_{k=0}^m \left(\sum_{i=0}^n p_{ki}^T p_{il} \right) \gamma_{lk}^{(-1)} \sigma^2 = \frac{1}{n+1} \sum_{l=0}^m \sum_{k=0}^m \gamma_{kl} \gamma_{lk}^{(-1)} \sigma^2 = \frac{m+1}{n+1} \sigma^2; \end{aligned}$$

здесь $\gamma_{kl} = \sum_{i=0}^n p_{ki}^T p_{il}$ — элементы матрицы Γ . ■

Замечание 1. Величина $M[\rho^2]$ представляет собой среднее по точкам x , значение дисперсии случайной ошибки $\Delta\Phi_m(x_i)$. Из теоремы 11.12 следует, что $M[\rho^2] \rightarrow 0$ при $n \rightarrow \infty$, т.е. при неограниченном росте числа наблюдений среднее значение дисперсии ошибки стремится к нулю.

Замечание 2. Равенство (11.111) подтверждает представление о том, что статистические свойства метода наименьших квадратов проявляются при $n \gg m$, т.е. тогда, когда число наблюдений много больше числа параметров модели.

Введем обозначения $\delta_m = \delta(\Phi_m^y, y)$, $\delta_m^0 = \delta(\Phi_m^{y^0}, y^0)$, $d_m = \delta(\Phi_m^y, y^0)$. Заметим, что величина d_m представляет собой среднеквадратичное уклонение многочлена наилучшего среднеквадратичного приближения Φ_m^y (построенного по вектору y исходных данных, содержащих случайные ошибки) от вектора y^0 точных значений функции f , которую мы и пытаемся оценить. Именно эту величину следует считать «истинной» мерой погрешности. Вычислим математические ожидания величин δ_m^2 и d_m^2 .

Теорема 11.13. Справедливы равенства

$$M[d_m^2] = (\delta_m^0)^2 + \frac{m+1}{n+1} \sigma^2, \quad (11.112)$$

$$M[\delta_m^2] = (\delta_m^0)^2 + \frac{n-m}{n+1} \sigma^2. \quad (11.113)$$

Доказательство. В силу леммы 11.1 имеет место равенство

$$s(\mathbf{a}, \mathbf{y}^0) = s(\mathbf{a}^0, \mathbf{y}^0) + \|\mathbf{P}\Delta\mathbf{a}\|_2^2,$$

которое в принятых выше обозначениях примет вид $d_m^2 = (\delta_m^0)^2 + \rho^2$. Вычислив математическое ожидание и использовав теорему 11.12, получим равенство (11.112).

Заметим теперь, что в силу леммы 11.1 справедливо равенство $s(\mathbf{a}^0, \mathbf{y}) = s(\mathbf{a}, \mathbf{y}) + \|\mathbf{P}\Delta\mathbf{a}\|_2^2$. Преобразуем левую часть этого равенства:

$$s(\mathbf{a}^0, \mathbf{y}) = \|(\mathbf{P}\mathbf{a}^0 - \mathbf{y}^0) + \boldsymbol{\varepsilon}\|_2^2 = \|\mathbf{P}\mathbf{a}^0 - \mathbf{y}^0\|_2^2 + 2(\mathbf{P}\mathbf{a}^0 - \mathbf{y}^0, \boldsymbol{\varepsilon}) + \|\boldsymbol{\varepsilon}\|_2^2.$$

Таким образом,

$$\delta_m^2 = \frac{1}{n+1} s(\mathbf{a}, \mathbf{y}) = (\delta_m^0)^2 + \frac{1}{n+1} \|\boldsymbol{\varepsilon}\|_2^2 - \rho^2 + \frac{2}{n+1} (\mathbf{P}\mathbf{a}^0 - \mathbf{y}^0, \boldsymbol{\varepsilon}).$$

Учитывая, что

$$M[\|\boldsymbol{\varepsilon}\|_2^2] = \sum_{i=0}^n M[\varepsilon_i^2] = (n+1)\sigma^2, \quad M[\rho^2] = \frac{m+1}{n+1} \sigma^2,$$

$$M[(\mathbf{P}\mathbf{a}^0 - \mathbf{y}^0, \boldsymbol{\varepsilon})] = (\mathbf{P}\mathbf{a}^0 - \mathbf{y}^0, M[\boldsymbol{\varepsilon}]) = 0,$$

получаем равенство (11.113). ■

5. О выборе степени обобщенного многочлена. Пусть функцию f можно аппроксимировать с достаточно высокой точностью ϵ обобщенным многочленом $\Phi_m(x) = \sum_{k=0}^m a_k \phi_k(x)$ некоторой степени m . Если эта степень заранее не известна, то возникает проблема выбора оптимальной степени аппроксимирующего многочлена Φ_m^y в условиях, когда исходные данные y , содержат случайные ошибки $\boldsymbol{\varepsilon}_r$.

Пусть $\phi_0, \phi_1, \dots, \phi_m$ — фиксированный набор базисных функций, линейно независимых в точках x_0, x_1, \dots, x_n . Предположим, что $n \gg m$, а ошибки $\boldsymbol{\varepsilon}_r$ удовлетворяют условиям (11.105), (11.106). Будем решать задачу наименьших квадратов для $m = 0, 1, 2, \dots$, постепенно увеличивая число параметров модели. Отметим, что значения среднеквадратичных уклонений δ_m и δ_m^0 должны с ростом m убывать. Действительно,

множество всех многочленов Φ_{m+1} степени $m+1$ включает в себя множество всех многочленов Φ_m степени m и поэтому

$$\delta_{m+1} = \min_{\Phi_{m+1}} \delta(\Phi_{m+1}, y) \leq \min_{\Phi_m} \delta(\Phi_m, y) = \delta_m, \quad m \geq 0;$$

$$\delta_{m+1}^0 = \min_{\Phi_{m+1}} \delta(\Phi_{m+1}, y^0) \leq \min_{\Phi_m} \delta(\Phi_m, y^0) = \delta_m^0, \quad m \geq 0.$$

Заметим также, что в силу теоремы 11.13 и закона больших чисел (при $n \gg 1$) справедливы следующие приближенные равенства:

$$d_m^2 = \delta(\Phi_m^y, y^0)^2 \approx (\delta_m^0)^2 + \frac{m+1}{n+1} \sigma^2, \quad (11.114)$$

$$\delta_m^2 = \delta(\Phi_m^y, y)^2 \approx (\delta_m^0)^2 + \frac{n-m}{n+1} \sigma^2. \quad (11.115)$$

Как нетрудно видеть, с ростом m первое слагаемое в правой части равенства (11.114) убывает, а второе возрастает. Поэтому следует ожидать, что величина d_m (именно она и характеризует уклонение вычисляемого многочлена Φ_m^y от функции f) с ростом m должна сначала убывать, а затем, достигнув своего минимума при некотором $m = m_0$, начнет возрастать.

Итак, существует оптимальная (в смысле критерия d_m) степень аппроксимации m_0 . Однако несмотря на то, что увеличивая m , можно получать все лучшее соответствие многочлена Φ_m^y с экспериментальными данными, следует иметь в виду, что начиная с $m = m_0$, многочлен будет все хуже соответствовать приближаемой функции.

Выбор оптимальной степени m_0 не представлял бы труда, если бы значения d_m можно было вычислять. Однако в действительности прямому вычислению поддаются только значения среднеквадратичного уклонения δ_m , анализируя которые мы и должны выбирать степень многочлена.

Предположим, что дисперсия σ^2 (или ее оценка) случайных ошибок ε_i известна. Пусть также известно, что для некоторого значения $m < n$ возможно приближение, для которого $\delta_m^0 \ll \sigma$. Тогда, начиная с этого значения m , согласно формуле (11.115) с учетом того, что $\frac{n-m}{n+1} \approx 1$, будет иметь место приближенное равенство

$$\delta_m \approx \sigma. \quad (11.116)$$

Таким образом, за оптимальное значение степени многочлена следует принять то значение m , при котором впервые будет выполнено приближенное равенство (11.116). При имеющейся информации лучший выбор вряд ли возможен.

Пусть теперь значение σ^2 дисперсии ошибок ε , нам не известно. Зато известно, что функция f либо представляет собой обобщенный многочлен Φ_m некоторой степени m_0 (и тогда $\delta_m^0 = c_0 = 0$ для всех $m \geq m_0$), либо же, начиная с некоторого $m = m_0$, увеличение степени многочлена в широком диапазоне значений $m_0 \leq m \leq m_1$ практически не влияет на качество приближения (т.е. $\delta_m^0 \approx c_0 = \text{const}$ для всех $m_0 \leq m \leq m_1$). Согласно формуле (11.115) справедливо приближенное равенство

$$\sigma_m^2 \equiv \frac{n+1}{n-m} \delta_m^2 \approx \frac{n+1}{n-m} (\delta_m^0)^2 + \sigma^2. \quad (11.117)$$

При $m \geq m_0$ величина, стоящая в правой части этого равенства, приближенно равна $\frac{n+1}{n-m} c_0^2 + \sigma^2$ и не убывает по m . Это наблюдение позволяет использовать в рассматриваемой ситуации следующее практическое правило выбора: за оптимальное значение степени многочлена следует принять то значение m , начиная с которого величина σ_m^2 стабилизируется или начинает возрастать.

Пример 11.16. Найдем оптимальную степень алгебраического многочлена для аппроксимации функции, заданной табл. 11.13.

Заметим, что фактически оптимальная степень $m = 1$ уже была найдена при решении примера 11.15. Основанием для такого вывода послужило сравнение среднеквадратичных уклонений $\delta_1 \approx 0.0486$ и $\delta_2 \approx 0.0481$ с оценкой уровня «шума» таблицы — грубый аналог критерия (11.116).

Попробуем получить тот же вывод, используя значения величин (11.117). Найдем дополнительно $P_0(x) = 0.425$ и $\delta_0 \approx 0.162$. Тогда $\sigma_0^2 = \frac{10}{9} \delta_0^2 \approx 0.292$, $\sigma_1^2 = \frac{10}{8} \delta_1^2 \approx 0.00295$, $\sigma_2^2 = \frac{10}{7} \delta_2^2 \approx 0.00381$. Так как $\sigma_0^2 > \sigma_1^2$ и $\sigma_1^2 < \sigma_2^2$, то за оптимальное значение степени следует принять $m = 1$. ▲

6. Нелинейная задача наименьших квадратов. Часто из физических или каких-либо других соображений следует, что зависимость $y = f(x)$ между величинами y и x должна хорошо описываться моделью

вида $y = g(x, a)$, где функция $g(x, a) = g(x, a_0, a_1, \dots, a_m)$ нелинейно зависит от параметров a_0, a_1, \dots, a_m .

Пусть функция $y = f(x)$ задана таблицей значений $y_i = f(x_i)$, $i = 0, 1, \dots, n$, где $n \gg m$. Тогда применение критерия наименьших квадратов приводит к задаче определения искомых параметров a_0, a_1, \dots, a_m из условия минимума функции

$$s(a, y) = \sum_{i=0}^n (g(x_i, a) - y_i)^2.$$

Нелинейная задача наименьших квадратов (особенно при большом числе параметров) весьма трудна для решения. Обычно для вычисления параметров a применяются специальные методы минимизации [40] (см. § 10.7).

В некоторых весьма специальных случаях решение нелинейной задачи наименьших квадратов можно свести к решению линейной задачи. Пусть, например, зависимость y от x ищется в виде $y = ae^{bx}$, где $a > 0$. Логарифмируя это равенство, приходим к линейной зависимости $\ln y = \ln a + bx$ величины $Y = \ln y$ от x . Теперь по таблице значений $Y_i = \ln y_i$, $i = 0, 1, \dots, n$, используя метод наименьших квадратов, можно легко определить значения $\ln a$ и b . Подчеркнем все же, что найденные таким образом значения параметров a и b отличаются от тех значений, которые могут быть найдены непосредственной минимизацией функции $\sum_{i=0}^n (a e^{bx_i} - y_i)^2$.

§ 11.14. Равномерное приближение функций

1. Задача о наилучшем равномерном приближении. Пусть $f(x)$ — заданная на отрезке $[a, b]$ непрерывная функция. Будем говорить, что многочлен $P_n(x)$ приближает функцию f равномерно на отрезке $[a, b]$ с точностью ε , если $\Delta(P_n) = \max_{[a, b]} |f(x) - P_n(x)| \leq \varepsilon$. Таким образом,

величина $\Delta(P_n)$ играет здесь роль погрешности приближения.

Естественно поставить следующую задачу среди всех многочленов фиксированной степени n найти многочлен $Q_n(x)$, для которого величина погрешности равномерного приближения минимальна, т.е. $\Delta(Q_n) \leq \Delta(P_n)$ для любого многочлена $P_n(x)$ степени n . Поставленная задача называется **задачей о наилучшем равномерном приближении**, а искомый многочлен $Q_n(x)$ — **многочленом наилучшего равномерного приближения**. Справедливо следующее утверждение.

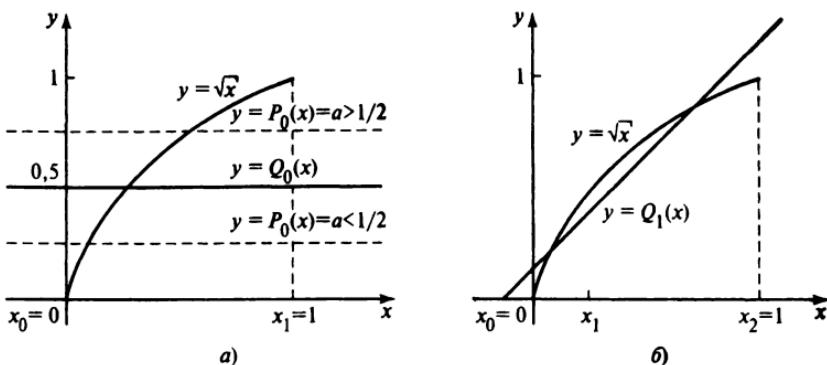


Рис. 11.12

Теорема 11.14. Для любой непрерывной на отрезке $[a, b]$ функции f многочлен наилучшего равномерного приближения Q_n степени n существует и единственен. ■

Пример 11.17. Покажем, что многочленом наилучшего равномерного приближения нулевой степени для функции $y = \sqrt{x}$ на отрезке $[0, 1]$ является $Q_0(x) = 1/2$.

Заметим, что $\Delta(Q_0) = \max_{[0, 1]} |\sqrt{x} - 1/2| = 1/2$ (рис. 11.12, а) и этот максимум достигается в точках $x_0 = 0$ и $x_1 = 1$. Для любого другого многочлена $P_0(x) = a$, где $a \neq 1/2$, значение $\Delta(P_0) > 1/2$. Действительно, если $a > 1/2$, то $\Delta(P_0) \geq |\sqrt{x_0} - a| = a > 1/2$. Если же $a < 1/2$, то $\Delta(P_0) \geq |\sqrt{x_1} - a| = 1 - a > 1/2$. ▲

Приведем без доказательства один из наиболее известных результатов о многочленах наилучшего равномерного приближения.

Теорема 11.15 (теорема Чебышева). Для того чтобы многочлен $Q_n(x)$ был многочленом наилучшего равномерного приближения непрерывной на отрезке $[a, b]$ функции $f(x)$, необходимо и достаточно, чтобы на отрезке $[a, b]$ нашлись по крайней мере $n + 2$ точки $x_0 < x_1 < x_2 < \dots < x_{n+1}$ такие, что

$$f(x_i) - Q_n(x_i) = \alpha(-1)^i \max_{[a, b]} |f(x) - Q_n(x)|, \quad i = 0, 1, \dots, n+1. \quad (11.118)$$

Здесь α — постоянная, равная 1 или -1 для всех i одновременно. ■

Точки x_0, x_1, \dots, x_{n+1} , удовлетворяющие условиям этой теоремы, принято называть *точками чебышевского альтернанса*¹. Равенство (11.118) накладывает на точки чебышевского альтернанса два требования:

1) в точках x_i модуль погрешности приближения функции f многочленом Q_n достигает максимума:

$$|f(x_i) - Q_n(x_i)| = \max_{[a, b]} |f(x) - Q_n(x)|;$$

2) для всех $i = 0, 1, \dots, n$ погрешность $f(x_i) - Q_n(x_i)$ меняет знак при переходе от точки x_i к следующей точке x_{i+1} .

Пример 11.18. В примере 11.17 точками чебышевского альтернанса являются точки $x_0 = 0$ и $x_1 = 1$ (число точек равно двум, так как $n = 0$). В самом деле, в этих точках достигается максимум модуля погрешности, а сама погрешность меняет знак при переходе от x_0 к x_1 . ▲

Пример 11.19. Найдем многочлен наилучшего равномерного приближения первой степени для функции $y = \sqrt{x}$ на отрезке $[0, 1]$.

В рассматриваемом случае $n = 1$ и должны быть по крайней мере три точки чебышевского альтернанса. В частности, это означает, что графики функций $y = \sqrt{x}$ и $y = Q_1(x) = ax + b$ должны пересекаться хотя бы дважды (рис. 11.12, б). Функция $\varphi(x) = \sqrt{x} - Q_1(x)$ вогнута и потому может иметь на отрезке $[0, 1]$ лишь одну внутреннюю точку экстремума x_1 . Следовательно, две из точек чебышевского альтернанса совпадают с концами отрезка: $x_0 = 0, x_2 = 1$. В точке x_1 функция $\varphi(x)$ удовлетворяет необходимому условию экстремума $\varphi'(x_1) = 0$, что эквивалентно уравнению

$$\frac{1}{2\sqrt{x_1}} - a = 0. \quad (11.119)$$

С учетом равенств $\varphi(x_0) = -b, \varphi(x_1) = \sqrt{x_1} - ax_1 - b, \varphi(x_2) = 1 - a - b$ условие перемены знака $\varphi(x_0) = -\varphi(x_1) = \varphi(x_2)$ эквивалентно двум уравнениям

$$-b = -\sqrt{x_1} + ax_1 + b, \quad (11.120)$$

$$-b = 1 - a - b. \quad (11.121)$$

¹ От латинского слова *alternare* — «чередоваться».

Из (11.121) сразу находим, что $a = 1$. Затем из (11.119) определяем, что $x_1 = 0.25$. Наконец, из (11.120) получаем $b = 1/8$.

Таким образом, $Q_1(x) = x + 1/8$ — многочлен наилучшего равномерного приближения первой степени, аппроксимирующий функцию \sqrt{x} с погрешностью $\Delta(Q_1) = 1/8$. ▲

2. Задача о понижении степени многочлена. Пусть $P_m(x) = a_m x^m + a_{m-1} x^{m-1} + \dots + a_0$ — многочлен степени $m \geq 1$, значения которого вычисляются на стандартном отрезке $[-1, 1]$ (напомним, что к стандартному отрезку можно перейти от произвольного отрезка $[a, b]$ линейной заменой переменных). Поставим следующую задачу о понижении степени многочлена: аппроксимировать $P_m(x)$ на отрезке $[-1, 1]$ многочленом $Q_{m-1}(x)$ наилучшего равномерного приближения на единицу меньшей степени.

Заметим, что $R_m(x) = P_m(x) - Q_{m-1}(x)$ — многочлен степени m со старшим коэффициентом, равным старшему коэффициенту a_m многочлена $P_m(x)$. Поставленную задачу можно иначе сформулировать так: среди всех многочленов степени m с фиксированным старшим коэффициентом a_m найти многочлен R_m , для которого величина $\max_{[-1, 1]} |R_m(x)|$ минимальна. Как известно (см. § 11.6), решение этой задачи дает многочлен $R_m(x) = \frac{a_m}{2^{m-1}} T_m(x)$ и для него $\max_{[-1, 1]} |R_m(x)| = \frac{|a_m|}{2^{m-1}}$. Таким образом, решением поставленной задачи является многочлен

$$Q_{m-1}(x) = P_m(x) - \frac{a_m}{2^{m-1}} T_m(x). \quad (11.122)$$

Соответствующая погрешность приближения равна $|a_m|/2^{m-1}$.

Замечание. Тривиальный способ понижения степени многочлена $P_m(x)$ — отбрасывание старшего слагаемого $a_m x^m$ — дает погрешность, равную $|a_m|$, т.е. в 2^{m-1} раз большую, чем приближение многочленом (11.122).

3. Нахождение многочленов, близких к наилучшим. В большинстве реальных случаев задача о наилучшем равномерном приближении непрерывной функции f является очень трудной. Для ее решения развиты специальные численные методы, реализованные в виде стандартных программ. Заметим, однако, что во многих ситуациях доста-

точно ограничиться нахождением многочлена, близкого к наилучшему, либо просто найти многочлен, равномерно приближающий функцию f с заданной точностью ε . Укажем на два случая, когда возможно нахождение многочленов, близких к наилучшим.

1°. Пусть производная $f^{(n+1)}(x)$ функции f слабо меняется на отрезке $[a, b]$. Тогда интерполяционный многочлен $P_n(x)$ с чебышевскими узлами (11.43) близок к многочлену наилучшего равномерного приближения.

2°. Пусть функция $f(x)$ задана на отрезке $[-1, 1]$ равномерно сходящимся степенным рядом (рядом Тейлора):

$$f(x) = \sum_{k=0}^{\infty} a_k x^k. \quad (11.123)$$

Требуется найти многочлен минимальной степени, равномерно приближающий функцию f на отрезке $[-1, 1]$ с заданной точностью ε .

Излагаемый ниже метод решения этой задачи часто называют *экономизацией степенных рядов*. Сначала берут отрезок ряда Тейлора

$$P_m(x) = \sum_{k=0}^m a_k x^k, \text{ аппроксимирующий функцию } f \text{ с точностью } \varepsilon_0 < \varepsilon.$$

Далее степень многочлена последовательно понижают. Если погрешность $\varepsilon_1 = |a_m|/2^{m-1}$ понижения степени такова, что $\varepsilon_0 + \varepsilon_1 \leq \varepsilon$, то многочлен P_m заменяют многочленом $P_{m-1} = P_m - \frac{a_m}{2^{m-1}} T_m(x) = \sum_{k=0}^{m-1} a_k^{(1)} x^k$.

Если погрешность $\varepsilon_2 = |a_{m-1}|/2^{m-2}$ такова, что $\varepsilon_0 + \varepsilon_1 + \varepsilon_2 \leq \varepsilon$, то снова понижают степень многочлена P_{m-1} по формуле $P_{m-2} = P_{m-1} - \frac{a_{m-1}^{(1)}}{2^{m-2}} T_{m-1}$ и т.д. Процесс прерывают тогда, когда вычисление очередного a_k дает $\varepsilon_0 + \varepsilon_1 + \dots + \varepsilon_k > \varepsilon$. В этом случае полагают $f(x) \approx P_{m-k+1}(x)$.

Пример 11.20. Найдем многочлен минимальной степени, аппроксимирующий функцию $y = \sin x$ на отрезке $[-1, 1]$ с точностью $\varepsilon = 10^{-3}$.

Ряд $\sin x = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}$ — знакопеременный, а его слагаемые убывают по модулю. Поэтому погрешность приближения функции $\sin x$ отрезком ряда Тейлора $P_{2n+1}(x) = \sum_{k=0}^n (-1)^k \frac{x^{2k+1}}{(2k+1)!}$ оценивается величиной $\varepsilon_0 = \frac{1}{(2n+3)!}$, равной максимуму модуля первого отброшенного

слагаемого. Выбор $n = 2$ дает значение $\varepsilon_0 = 1/5040 \approx 2 \cdot 10^{-4} \leq \varepsilon$. Следовательно, многочлен $P_5(x) = x - \frac{x^3}{6} + \frac{x^5}{120}$ аппроксимирует функцию $y = \sin x$ с точностью $\varepsilon_0 \approx 2 \cdot 10^{-4}$.

Понижение степени многочлена P_5 будет сопровождаться дополнительной погрешностью $\varepsilon_1 = \frac{|a_5|}{2^4} = \frac{1}{1920} \approx 5.2 \cdot 10^{-4}$. Следовательно, $\varepsilon_0 + \varepsilon_1 < \varepsilon$ и после понижения степени по формуле

$$P_5(x) - \frac{1}{1920} T_5(x) = x - \frac{x^3}{6} + \frac{x^5}{120} - \frac{1}{1920} (16x^5 - 20x^3 + 5x)$$

получим многочлен

$$P_3(x) = \frac{383}{384} x - \frac{5}{32} x^3, \quad (11.124)$$

дающий приближение к $\sin x$ на отрезке $[-1, 1]$ с точностью $7.2 \cdot 10^{-4}$.

Так как $\varepsilon_3 = \frac{5}{32} \cdot \frac{1}{2^2} > 10^{-3}$, то дальнейшее понижение степени невозможно и решением задачи является многочлен (11.124). ▲

§ 11.15. Дробно-рациональные аппроксимации и вычисление элементарных функций

При создании стандартных программ вычисления элементарных и специальных функций специалистами по математическому обеспечению компьютеров применяются разнообразные приемы, требующие глубоких профессиональных знаний. Используемые вычислительные методы не являются здесь машинно-независимыми, а, наоборот, существенно учитывают разрядность мантиссы, скорость выполнения арифметических операций и другие особенности конкретного компьютера. Отметим, что к указанным стандартным программам обычно предъявляется требование обеспечения относительной точности результата порядка машинного эпсилон ε_m .

Использование богатой дополнительной информации об аналитических свойствах элементарных и специальных функций позволяет значительно уменьшить объем вычислений. Существенно используется возможность представления вычисляемых функций сходящимися степенными рядами вида

$$\sum_{k=0}^{\infty} c_k z^k. \quad (11.125)$$

Здесь $z = x - x_0$; x_0 — точка, в которой осуществляется разложение функции в ряд. Отметим, однако, что вопреки распространенному мнению такие ряды непосредственно практически никогда не используются для вычисления функций.

Широко применяемым в настоящее время способом представления функций является приближение их рациональными дробями вида

$$R(z) = \frac{a_0 + a_1 z + \dots + a_n z^n}{b_0 + b_1 z + \dots + b_m z^m}. \quad (11.126)$$

Кдробно-рациональным аппроксимациям приходят различными путями. В ряде случаев используется *рациональная интерполяция* — интерполяция функции рациональной дробью (11.126). Тогда коэффициенты a_j ($j = 0, 1, \dots, n$), b_k ($k = 0, 1, \dots, m$) находятся из совокупности соотношений $R(z_i) = y_i$ ($0 \leq i < N$, $N = n + m + 1$), которые можно записать в следующем виде:

$$\sum_{j=0}^n a_j x_i^j = y_i \sum_{k=0}^m b_k x_i^k, \quad 0 \leq i < N.$$

Эти соотношения образуют систему N линейных однородных алгебраических уравнений относительно $N+1$ неизвестных — коэффициентов $a_0, \dots, a_n, b_0, \dots, b_m$. Такие системы всегда имеют нетривиальные решения. И хотя таких решений бесконечно много, найденная с их помощью рациональная дробь определяется однозначно. Чтобы избежать неоднозначности в определении коэффициентов этой дроби достаточно использовать условие нормировки $b_0 = 1$. Можно записать $R(z)$ в явном виде, если использовать аппарат разделенных разностей [9].

Один из возможных путей состоит в использовании теории *цепных* (или непрерывных) дробей. Например, функция $\operatorname{tg} x$ представляется цепной дробью

$$\operatorname{tg} x = \cfrac{x}{1 - \cfrac{x^2}{3 - \cfrac{x^2}{5 - \cfrac{\ddots}{}}}}.$$

Обрывая такую бесконечную дробь, получают некоторую конечную дробь, аппроксимирующую функцию.

Все более популярным в последние годы способом приближения аналитических функций становится *аппроксимация Паде* — такая дробно-линейная аппроксимация (11.126), для которой

$$\sum_{k=0}^{\infty} c_k z^k = \frac{a_0 + a_1 z + \dots + a_n z^n}{b_0 + b_1 z + \dots + b_m z^m} + O(z^{n+m+1}). \quad (11.127)$$

Равенство (11.127) означает, что коэффициенты $a_0, a_1, \dots, a_n, b_0, b_1, \dots, b_m$ дроби (11.126) подбираются так, чтобы в разложении ее в ряд Тейлора первые $n + m + 1$ слагаемых в точности совпадали с соответствующими слагаемыми ряда (11.125).

Например, аппроксимация Паде функции e^x в случае $n = m \geq 1$ задается формулой

$$e^x \approx R_n(x) = \frac{P_n(x)}{P_n(-x)},$$

где

$$P_n(x) = 1 + \frac{1}{2} x + \frac{n-1}{2(2n-1)} \frac{x^2}{2!} + \dots + \frac{(n-1) \cdot \dots \cdot 1}{2(2n-1) \cdot \dots \cdot (n+1)} \frac{x^n}{n!},$$

с погрешностью

$$e^x - R_n(x) = (-1)^n \frac{(n!)^2}{(2n)!(2n+1)!} x^{2n+1} + O(x^{2n+2}).$$

При $n = m = 2$ и $n = m = 3$ имеем:

$$\begin{aligned} e^x &\approx \frac{12 + 6x + x^2}{12 - 6x + x^2}, \\ e^x &\approx \frac{12(x^2 + 10) + x(x^2 + 60)}{12(x^2 + 10) - x(x^2 + 60)}. \end{aligned} \quad (11.128)$$

Отметим, что при $|x| \leq 0.5 \ln 2$ аппроксимация (11.128) обеспечивает точность $9 \cdot 10^{-9}$.

В заключение параграфа рассмотрим, как вычисляется в одной из стандартных программ функция $y = \ln x$. Аргумент x представляют в виде $x = \mu \cdot 2^p$, где p — двоичный порядок, μ — мантисса, $0.5 \leq \mu < 1$. Затем используют разложение

$$\ln x = (p - 0.5) \ln 2 + 2v \sum_{k=0}^{\infty} \frac{1}{2k+1} v^{2k},$$

где $v = (\sqrt{2}\mu - 1)/(\sqrt{2}\mu + 1)$.

Заметим, что $|v| \leq a = (\sqrt{2} - 1)/(\sqrt{2} + 1) < 0.172$. Ряд $\sum_{k=0}^{\infty} \frac{1}{2k+1} v^{2k}$

заменяют затем на отрезке $[-a, a]$ многочленом, близким к многочлену четвертой степени наилучшего равномерного приближения. В результате получается формула

$\ln x \approx (p - 0.5) \ln 2 + v(2.000000815 + 0.666445069v^2 + 0.415054254v^4)$,
погрешность которой для всех x не превышает $3 \cdot 10^{-8}$.

§ 11.16. Дополнительные сведения об интерполяции

1. Обратная интерполяция. Пусть функция $y = f(x)$ обратима. Тогда всякую таблицу ее значений $y_i = f(x_i)$, $i = 0, 1, \dots, n$ можно интерпретировать как таблицу значений $x_i = f^{-1}(y_i)$, $i = 0, 1, \dots, n$ обратной функции $x = f^{-1}(y)$. *Обратной интерполяцией* принято называть построение функции $x = p(y)$, интерполирующей обратную функцию, т.е. удовлетворяющей условию

$$p(y_i) = x_i, \quad i = 0, 1, \dots, n.$$

Этот подход часто используется для приближенного решения уравнения $f(x) = 0$. Достаточно построить таблицу значений функции f на отрезке локализации, выполнить обратную интерполяцию и положить $x \approx p(0)$.

2. Интерполяционный многочлен Бесселя. Иногда точность интерполяции можно повысить, если использовать линейные комбинации интерполяционных многочленов.

Пусть, например, функция f задана таблицей своих значений $y_i = f(x_i)$, $i = 0, \dots, n$ с постоянным шагом h и требуется найти приближенное значение функции в точке $x \in [x_{j-1}, x_j]$, используя интерполяцию многочленом второй степени. Естественно воспользоваться многочленом $P_{(j-2, j-1, j)}(x)$ с узлами интерполяции x_{j-2}, x_{j-1}, x_j или многочленом $P_{(j-1, j, j+1)}(x)$ с узлами x_{j-1}, x_j, x_{j+1} соответственно. Заметим, что каждый из этих многочленов использует узлы, которые расположены асимметрично относительно отрезка $[x_{j-1}, x_j]$. При этом для погрешности интерполяции справедливы оценки

$$\max_{[x_{j-1}, x_j]} |P_{(j-2, j-1, j)}(x) - f(x)| \leq \frac{M_3}{16} h^3,$$

$$\max_{[x_{j-1}, x_j]} |P_{(j-1, j, j+1)}(x) - f(x)| \leq \frac{M_3}{16} h^3.$$

Как обычно, $M_n = \max |f^{(n)}(x)|$.

Построим теперь *интерполяционный многочлен Бесселя* второй степени как полусумму $B(x) = \frac{1}{2}(P_{(j-2, j-1, j)}(x) + P_{(j-1, j, j+1)}(x))$. Стого говоря, этот многочлен не является интерполяционным, так он как совпадает с функцией f лишь в двух точках: x_{j-1} и x_j . Однако благодаря симметричному расположению использованных узлов $x_{j-2}, x_{j-1}, x_j, x_{j+1}$ относительно отрезка $[x_{j-1}, x_j]$ точность приближения оказывается выше:

$$\max_{[x_{j-1}, x_j]} |B(x) - f(x)| \leq \frac{M_3}{72\sqrt{3}} h^3 + \frac{3M_4}{128} h^4.$$

Аналогичным образом можно использовать интерполяционные многочлены Бесселя более высоких четных степеней.

3. Многомерная интерполяция. Интерполяция широко применяется для приближения функций многих переменных. Дадим понятие о многомерной интерполяции на примере функции двух переменных.

Предположим, что функция $z = f(x, y)$ задана на прямоугольнике $[a, b] \times [c, d]$ таблицей значений

$$z_{ij} = f(x_i, y_j), \quad 0 \leq i \leq n, \quad 0 \leq j \leq m.$$

Интерполяционный многочлен $P(x, y)$, обладающий свойством

$$P(x_i, y_j) = z_{ij}, \quad 0 \leq i \leq n, \quad 0 \leq j \leq m,$$

может быть записан в следующем виде:

$$L_{n, m}(x, y) = \sum_{i=0}^n \sum_{j=0}^m z_{ij} l_{ni}(x) l_{mj}(y). \quad (11.129)$$

$$\text{Здесь } l_{ni}(x) = \prod_{\substack{s=0 \\ s \neq i}}^n \frac{x - x_s}{x_i - x_s}, \quad l_{mj}(y) = \prod_{\substack{k=0 \\ k \neq j}}^m \frac{y - y_k}{y_j - y_k}.$$

Как нетрудно видеть, многочлен (11.129) является аналогом многочлена Лагранжа (11.22), причем при фиксированном значении y он является многочленом степени n по переменной x , а при фиксированном значении x — многочленом степени m по переменной y .

При $n = m = 1$ этот многочлен называется *билинейным*, так как он линеен по каждой из переменных. Приведем явную формулу, соответствующую *билинейной интерполяции*.

$$\begin{aligned} L_{1,1}(x, y) = & \frac{1}{(x_1 - x_0)(y_1 - y_0)} [z_{0,0}(x_1 - x)(y_1 - y) + z_{1,0}(x - x_0)(y_1 - y) + \\ & + z_{0,1}(x_1 - x)(y - y_0) + z_{1,1}(x - x_0)(y - y_0)]. \end{aligned}$$

Одна из проблем, связанных с использованием интерполяционного многочлена (11.129), состоит в предположении, что функция f задана на прямоугольнике, а узлы интерполяции образуют правильную прямоугольную сетку. Однако часто область определения функции имеет более сложную геометрическую форму, а узлы, в которых известны значения функции f , расположены по иному закону. В такой ситуации вполне приемлемым способом приближения функции f в точке (x, y) может оказаться *линейная интерполяция* по значениям в ближайших трех узлах. Пусть значения функции f известны в вершинах $A = (x_A, y_A)$, $B = (x_B, y_B)$ и $C = (x_C, y_C)$ треугольника Δ :

$$z_A = f(x_A, y_A), \quad z_B = f(x_B, y_B), \quad z_C = f(x_C, y_C).$$

Соответствующий интерполяционный многочлен имеет вид

$$P_1(x, y) = z_A + \alpha(x - x_A) + \beta(y - y_A), \quad (11.130)$$

где

$$\alpha = -[(z_C - z_A)(y_B - y_A) - (z_B - z_A)(y_C - y_A)]/\gamma,$$

$$\beta = [(x_B - x_A)(z_C - z_A) - (x_C - x_A)(z_B - z_A)]/\gamma,$$

$$\gamma = (x_B - x_A)(y_C - y_A) - (x_C - x_A)(y_B - y_A).$$

Графиком многочлена P_1 является плоскость, проходящая через точки с координатами (x_A, y_A, z_A) , (x_B, y_B, z_B) и (x_C, y_C, z_C) . Из курса аналитической геометрии известно, что уравнение такой плоскости может быть записано в виде

$$\begin{vmatrix} x - x_A & y - y_A & z - z_A \\ x_B - x_A & y_B - y_A & z_B - z_A \\ x_C - x_A & y_C - y_A & z_C - z_A \end{vmatrix} = 0.$$

Как нетрудно видеть, это уравнение эквивалентно формуле (11.30).

§ 11.17. Дополнительные замечания

1. Тем, кто интересуется практическим применением сплайнов, рекомендуем книгу [18]. В ней можно найти не только изложение теории и алгоритмов решения задач, но и полезные практические советы и тексты соответствующих программ на языке Фортран.

2. Изложение современных численных методов решения линейных задач метода наименьших квадратов содержится в общепризнанном учебнике-справочнике [64]. Заметим, что в оригинальное издание этой книги были включены и тексты соответствующих программ, которые были к сожалению, опущены при переводе.

3. Исключительно богатый материал по специальным и элементарным функциям содержит справочник [92]. Эта книга пользуется заслуженной популярностью у специалистов, и тот, кто занимается вычислением специальных функций лишь эпизодически, может найти в ней ответы на большинство интересующих его вопросов. Правда, большая часть приведенных в справочнике результатов была получена до 1960 г. Как дополнение к [92] можно рассматривать книгу [65].

4. Желающим более глубоко ознакомиться с аппроксимациями Паде и непрерывными (цепными) дробями рекомендуем обратиться к специальной литературе [11, 39].

5. В машинном проектировании широкое распространение получил подход к описанию кривых и поверхностей, заданных графически, с помощью *кривых Безье*. Этот подход были разработан в 1962 г. для автомобильных фирм Рено и Ситроен с целью упростить дизайнеру интерактивный режим работы. Первоначальное представление о кривых Безье можно получить из [54, 122].

6. В последние два десятилетия бурное развитие получила теория вэйвлетов (всплесков). Вэйвлет-анализ нашел широкое применение в цифровой обработке информации и изображений, в теории фильтрации и кодирования. Книга [113] содержит не только основы теории вейвлетов, но и хорошую библиографию по этой теме; см. также [8*, 13*].

7. Подчеркнем, что в данной главе рассматривались методы приближения функций только одной вещественной переменной (за исключением § 11.16). При аппроксимации функций многих переменных можно использовать аналогичные подходы (см., например, [9, 67]).

ЧИСЛЕННОЕ ДИФФЕРЕНЦИРОВАНИЕ

Численное дифференцирование применяется тогда, когда функцию трудно или невозможно продифференцировать аналитически. Например, необходимость в численном дифференцировании возникает в том случае, когда функция задана таблицей. Кроме того, формулы численного дифференцирования широко используются при разработке вычислительных методов решения многих задач (решение дифференциальных уравнений, поиск решений нелинейных уравнений, поиск точек экстремума функций и др.).

§ 12.1. Простейшие формулы численного дифференцирования

1. Вычисление первой производной. Предположим, что в окрестности точки x функция f дифференцируема достаточное число раз. Исходя из определения первой производной

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x},$$

естественно попытаться использовать для ее вычисления две простейшие приближенные формулы

$$f'(x) \approx \frac{f(x + h) - f(x)}{h}, \quad (12.1)$$

$$f'(x) \approx \frac{f(x) - f(x - h)}{h}, \quad (12.2)$$

соответствующие выбору фиксированных значений $\Delta x = h$ и $\Delta x = -h$. Здесь $h > 0$ — малый параметр (*шаг*). Разностные отношения в правых частях формул (12.1) и (12.2) часто называют *правой* и *левой разностными производными* [85].

Для оценки погрешностей

$$r_+(x, h) = f'(x) - \frac{f(x + h) - f(x)}{h},$$

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

введенных формул численного дифференцирования (*погрешностей аппроксимации*) воспользуемся формулами Тейлора.

$$f(x \pm h) = f(x) \pm f'(x)h + \frac{f''(\xi_{\pm})}{2} h^2. \quad (12.3)$$

Здесь и ниже ξ_+ и ξ_- — некоторые точки, расположенные на интервалах $(x, x+h)$ и $(x-h, x)$ соответственно. Подставляя разложения (12.3) в выражения для r_{\pm} , получаем $r_+(x, h) = -\frac{1}{2} f''(\xi_+)h$, $r_-(x, h) = \frac{1}{2} f''(\xi_-)h$.

Следовательно,

$$|r_+(x, h)| \leq \frac{1}{2} M_2 h, \quad M_2 = \max_{[x, x+h]} |f''(\xi)|, \quad (12.4)$$

$$|r_-(x, h)| \leq \frac{1}{2} M_2 h, \quad M_2 = \max_{[x-h, x]} |f''(\xi)|. \quad (12.5)$$

Таким образом, формулы (12.1), (12.2) имеют первый порядок точности по h . Иначе говоря, правая и левая разностные производные аппроксимируют производную $f'(x)$ с первым порядком точности.

Приведенные формулы численного дифференцирования имеют простую геометрическую интерпретацию (рис. 12.1, а). Пусть N_0 , N_- и N_+ — расположенные на графике функции $y = f(x)$ точки с координатами $(x, f(x))$, $(x-h, f(x-h))$ и $(x+h, f(x+h))$. Напомним, что производная $f'(x)$ равна тангенсу угла α наклона к оси Ox касательной, проведенной к графику функции в точке N_0 . Формула (12.1) соответствует приближенной замене производной $f'(x) = \operatorname{tg} \alpha$ правой разностной производной $\frac{f(x+h) - f(x)}{h}$, равной тангенсу угла α_+ наклона к графику функции секущей, проведенной через точки N_0 и N_+ . Формула (12.2) соответствует аналогичной замене левой разностной производной $\frac{f(x) - f(x-h)}{h}$, равной тангенсу угла α_- секущей, проведенной через точки N_0 и N_- .

Естественно предположить (рис. 12.1, а и б), что лучшим по сравнению с $\operatorname{tg} \alpha_+$ и $\operatorname{tg} \alpha_-$ приближением к $f'(x) = \operatorname{tg} \alpha$ является тангенс угла наклона α_0 секущей к графику, проведенной через точки N_- и N_+ . Соответствующая приближенная формула имеет вид

$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}. \quad (12.6)$$

Величину в правой части этой формулы часто называют *центральной разностной производной*.

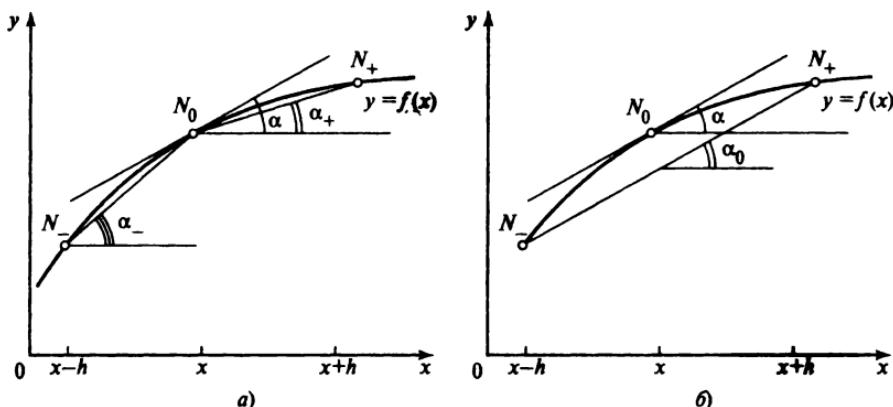


Рис. 12.1

Подставляя в выражение для погрешности

$$r_0(x, h) = f'(x) - \frac{f(x+h) - f(x-h)}{2h}$$

соответствующие разложения по формуле Тейлора

$$f(x \pm h) = f(x) \pm f'(x)h + \frac{f''(x)}{2} h^2 \pm \frac{f^{(3)}(\xi_{\pm})}{6} h^3,$$

получаем

$$r_0(x, h) = -\frac{f^{(3)}(\xi_+) + f^{(3)}(\xi_-)}{12} h^2.$$

Следовательно, справедлива оценка погрешности

$$|r_0(x, h)| \leq \frac{M_3}{6} h^2, \quad M_3 = \max_{[x-h, x+h]} |f^{(3)}(\xi)|. \quad (12.7)$$

Таким образом, центральная разностная производная аппроксимирует производную $f'(x)$ со вторым порядком точности относительно h .

Для вычисления $f'(x)$ можно получить формулы любого порядка точности (см. § 12.2). Однако в таких формулах с ростом порядка точности возрастает и число используемых значений функции. В качестве примера приведем формулу

$$f'(x) \approx \frac{f(x-2h) - 8f(x-h) + 8f(x+h) - f(x+2h)}{12h}, \quad (12.8)$$

имеющую четвертый порядок точности.

Пример 12.1. Пусть функция $f(x) = e^x$ задана на отрезке $[0, 1]$ таблицей значений с шагом $h = 0.2$ (табл. 12.1).

Таблица 12.1

x	0.0	0.2	0.4	0.6	0.8	1.0
$f(x)$	1.000000	1.22140	1.49182	1.82212	2.22554	2.71828

Используя формулы численного дифференцирования, найдем значения производной $f'(x)$ в узлах таблицы.

В точках $x = 0.0$ и $x = 1.0$ из приведенных в этом параграфе формул можно воспользоваться лишь формулами (12.1) и (12.2). В остальных точках применим формулу (12.6), имеющую более высокий порядок точности. Вычисления дают следующую таблицу производных (табл. 12.2).

Таблица 12.2

x	0.0	0.2	0.4	0.6	0.8	1.0
$f'(x)$	1.10700	1.22955	1.50180	1.83430	2.24040	2.46370
$r(x)$	-0.10700	-0.00815	-0.00998	-0.01218	-0.01486	0.25458

Здесь же приведены значения погрешностей $r(x)$, которые в данном случае легко вычисляются (ведь $f'(x) = e^x = f(x)$). Как и следовало ожидать, значения погрешности в крайних точках (здесь использовались формулы первого порядка точности) существенно больше, чем во внутренних точках.

Заметим, что значения погрешностей можно было оценить и заранее, используя априорные оценки (12.4), (12.5), (12.7). Например, неравенство (12.4) дает в точке $x = 0$ оценку

$$|r| \leq \frac{1}{2} M_2 h = \frac{1}{2} \max_{[0, 0.2]} |e^\xi| \cdot h = \frac{1}{2} \cdot 1.22140 \cdot 0.2 \approx 0.13 .$$

В точке $x = 0.2$ из неравенства (12.7) следует, что

$$|r| \leq \frac{1}{6} M_3 h^2 = \frac{1}{6} \max_{[0, 0.4]} |e^\xi| \cdot h^2 = \frac{1}{6} \cdot 1.49182 \cdot 0.04 \approx 0.01 , \text{ и т.д.}$$

Для вычисления $f'(x)$ при $x = 0.4$ и $x = 0.6$ можно также применить формулу (12.8) и получить значения $f'(x) \approx 1.49176$ и $f'(0.6) \approx 1.82203$ с погрешностями, приближенно равными $6 \cdot 10^{-5}$ и $9 \cdot 10^{-5}$ соответственно. ▲

2. Вычисление второй производной. Наиболее простой и широко применяемой для приближенного вычисления второй производной является следующая формула:

$$f''(x) \approx \frac{f(x-h) - 2f(x) + f(x+h)}{h^2}. \quad (12.9)$$

Величину в правой части этого приближенного равенства часто называют *второй разностной производной* [85].

Подставляя в выражение для погрешности

$$r(x, h) = f''(x) - \frac{f(x-h) - 2f(x) + f(x+h)}{h^2}$$

соответствующие разложения по формуле Тейлора

$$f(x \pm h) = f(x) \pm f'(x)h + \frac{f''(x)}{2}h^2 \pm \frac{f^{(3)}(x)}{6}h^3 + \frac{f^{(4)}(\xi_{\pm})}{24}h^4,$$

получаем

$$r(x, h) = -\frac{f^{(4)}(\xi_+) + f^{(4)}(\xi_-)}{24}h^2.$$

Следовательно,

$$|r(x, h)| \leq \frac{M_4}{12}h^2, \quad M_4 = \max_{[x-h, x+h]} |f^{(4)}(\xi)|. \quad (12.10)$$

Таким образом, формула (12.9) имеет второй порядок точности.

Для вычисления $f''(x)$ можно получить формулы любого порядка точности. Например, формула

$$f''(x) \approx \frac{-f(x-2h) + 16f(x-h) - 30f(x) + 16f(x+h) - f(x+2h)}{12h^2} \quad (12.11)$$

имеет четвертый порядок точности.

Пример 12.2. Используя табл. 12.1 значений функции $f(x) = e^x$, найдем с помощью формул численного дифференцирования значения $f''(x)$ во внутренних узлах таблицы.

Вычисление по формуле (12.9) дает значения, приведенные в табл. 12.3.

Таблица 12.3

x	0.2	0.4	0.6	0.8
$f''(x)$	1.22550	1.49700	1.82800	2.23300
$r(x)$	-0.00410	-0.00518	-0.00588	-0.00746

Применение формулы (12.11) позволяет получить значения $f''(0.4) \approx 1.49204$, $f''(0.6) \approx 1.82183$ с погрешностями, приближенно равными $-2 \cdot 10^{-4}$, $-3 \cdot 10^{-4}$. ▲

§ 12.2. О выводе формул численного дифференцирования

Хотя простейшие формулы численного дифференцирования можно получить сравнительно элементарно, для вывода и анализа таких формул в более сложных случаях необходимо использовать значительно более серьезный математический аппарат. Заметим, что основой для построения различных приближенных формул вычисления производных являются методы теории приближения функций, элементы которой были изложены в предыдущей главе.

Предположим, что в окрестности точки x функция аппроксимируется некоторой другой функцией g , причем производная $g^{(k)}$ в точке x легко вычисляется. Естественно в такой ситуации попытаться воспользоваться приближенной формулой

$$f^{(k)}(x) \approx g^{(k)}(x). \quad (12.12)$$

Наиболее просто этот подход реализуется, когда приближение осуществляется с помощью интерполяции.

1. Формулы численного дифференцирования, основанные на интерполяции алгебраическими многочленами. Пусть $P_n(x)$ — интерполяционный многочлен степени n с узлами интерполяции $x_0 < x_1 < x_2 < \dots < x_n$ и $x \in [x_0, x_n]$. В этом случае формула (12.12) принимает вид

$$f^{(k)}(x) \approx P_n^{(k)}(x), \quad 0 \leq k \leq n. \quad (12.13)$$

При этом справедлива следующая оценка погрешности формулы (12.13):

$$|f^{(k)}(x) - P_n^{(k)}(x)| \leq C_{n,k} M_{n+1} h_{\max}^{n+1-k}, \quad 0 \leq k \leq n, \quad (12.14)$$

где $C_{n,k}$ — положительные числа, $M_{n+1} = \max_{[x_0, x_n]} |f^{n+1}(x)|$.

Замечание 1. Порядок точности формулы (12.13) относительно h_{\max} равен разности между числом узлов интерполяции и порядком вычисляемой производной.

Замечание 2. Если формула (12.13) применяется для вычисления производной в точке, относительно которой узлы таблицы расположены симметрично, и число $n - k$ четно, то порядок точности формулы повышается на единицу по сравнению с порядком $n + 1 - k$, гарантированным оценкой (12.14). Таковы, например формулы (12.6), (12.8), (12.9), (12.11).

Заметим, что $P_n^{(n)}(x) = n! f(x_0; x_1; \dots, x_n)$ (это следует, в частности, из формулы (11.52)). Таким образом, справедлива приближенная формула (она вытекает и из (11.50))

$$f^{(n)}(x) \approx n! f(x_0; x_1; \dots, x_n), \quad (12.15)$$

имеющая по крайней мере первый порядок точности. Ее частными случаями являются следующие формулы:

$$f'(x) \approx f(x_0; x_1) = \frac{f(x_1) - f(x_0)}{x_1 - x_0}, \quad (12.16)$$

$$f''(x) \approx 2f(x_0; x_1; x_2) = \frac{2}{x_2 - x_0} \left[\frac{f(x_2) - f(x_1)}{x_2 - x_1} - \frac{f(x_1) - f(x_0)}{x_1 - x_0} \right]. \quad (12.17)$$

При выборе в качестве узлов интерполяции значений $x_0 = x$, $x_1 = x + h$ формула (12.16) превращается в формулу (12.1). При выборе $x_0 = x - h$, $x_1 = x$ из (12.16) получается формула (12.2), а при $x_0 = x - h$, $x_1 = x + h$ — формула (12.6). Аналогично из формулы (12.17) получается формула (12.9).

2. Использование таблиц с постоянным шагом. Наиболее простой вид принимают формулы численного дифференцирования при использовании таблиц $y_i = f(x_i)$ с постоянным шагом. Например, формула (12.15) в этом случае выглядит так:

$$f^{(n)}(x) \approx \frac{\Delta^n y_0}{h^n}.$$

В тех случаях, когда значение производной необходимо вычислять в крайних для таблицы точках x_0 и x_n , используются *односторонние формулы численного дифференцирования* $f^{(k)}(x_0) \approx P_n^{(k)}(x_0)$ и $f^{(k)}(x_n) \approx P_n^{(k)}(x_n)$. Приведем односторонние формулы (их легко получить дифференцированием многочленов Ньютона (11.57) и (11.58)) для вычисления первой производной f' :

$$f'(x_0) \approx \frac{1}{h} \sum_{j=1}^n \frac{(-1)^{j-1}}{j} \Delta^j y_0, \quad (12.18)$$

$$f'(x_n) \approx \frac{1}{h} \sum_{j=1}^n \frac{1}{j} \nabla^j y_n, \quad (12.19)$$

имеющие n -й порядок точности.

При $n = 2$ из (12.18), (12.19) получаются формулы

$$f'(x_0) \approx \frac{1}{2h} (-3f(x_0) + 4f(x_1) - f(x_2)), \quad (12.20)$$

$$f'(x_n) \approx \frac{1}{2h} (f(x_{n-2}) - 4f(x_{n-1}) + 3f(x_n)), \quad (12.21)$$

имеющие второй порядок точности.

Пример 12.3. Вычисление значений производной функции $f(x) = e^x$ заданной табл. 12.1, по формулам (12.20), (12.21) дает значения $f'(0.0) \approx 0.98445$, $f'(1.0) \approx 2.68700$ с погрешностями, равными 0.01555 и 0.03128. Для сравнения напомним, что погрешности значений, найденных в примере 12.1 по простейшим формулам (12.1), (12.2) соответственно равны — 0.10700 и 0.25458. ▲

3. Другие подходы. Применение формулы (12.13) для вычисления производной $f^{(k)}$ фактически основано на кусочно-полиномиальной интерполяции. Полученная таким образом производная в точке «стыка» двух соседних многочленов может иметь разрывы. Поэтому, если требуется глобально на отрезке $[a, b]$ аппроксимировать производную гладкой функцией, то целесообразно использовать сплайны. Производная $S_m^{(k)}(x)$ сплайна $S_m(x)$ при $k \leq m - r$ (где r — дефект сплайна) дает гладкую глобальную аппроксимацию для $f^{(k)}(x)$.

Когда значения функции сильно «зашумлены» случайными ошибками, полезным может оказаться использование метода наименьших квадратов.

§ 12.3. Обусловленность формул численного дифференцирования

Несмотря на внешнюю простоту формул численного дифференцирования, их применение требует особой осторожности. Отметим, что используемые при численном дифференцировании значения $f^*(x)$ функции $f(x)$ непременно содержат погрешности. Поэтому к погрешности аппроксимации формул численного дифференцирования добавляется неустранимая погрешность, вызванная погрешностями вычисления функции f . Для того чтобы погрешность аппроксимации была достаточно малой, требуется использование таблиц с малыми шагами h . Однако, к сожалению, при малых шагах формулы численного дифференцирования становятся плохо обусловленными и результат их применения может быть полностью искажен неустранимой погрешностью. Важно понимать, что действительная причина этого явления лежит не

в несовершенстве предложенных методов вычисления производных, а в некорректности самой операции дифференцирования приближенно заданной функции (см. гл. 3, пример 3.5).

Поясним сказанное на примере использования формулы (12.1). Полная погрешность $r^*(x, h) = f'(x) - \frac{f^*(x+h) - f^*(x)}{h}$ реально вычисляемого значения правой разностной производной представляет собой сумму погрешности аппроксимации $r_+(x, h) = f'(x) - \frac{f(x+h) - f(x)}{h}$ и неустранимой погрешности $r_h(x, h) = \frac{1}{h} [(f(x+h) - f^*(x+h)) - (f(x) - f^*(x))]$.

Пусть $\bar{\Delta}$ — верхняя граница абсолютной погрешности $\Delta(f^*(x)) = |f(x) - f^*(x)|$ используемых значений функции. Тогда неустранимая погрешность r_h оценивается следующим образом:

$$|r_h| \leq 2\bar{\Delta}/h. \quad (12.22)$$

Оценка (12.22) означает, что чувствительность формулы (12.1) к погрешностям входных данных характеризуется абсолютным числом обусловленности $v_\Delta = 2/h$. Так как $v_\Delta \rightarrow \infty$ при $h \rightarrow 0$, то формула (12.1) при малых h становится очень плохо обусловленной. Поэтому несмотря на то, что погрешность аппроксимации стремится к нулю при $h \rightarrow 0$ (см. оценку (12.4)), следует ожидать, что полная погрешность будет неограниченно возрастать при $h \rightarrow 0$. Во всяком случае так ведет себя верхняя граница полной погрешности $\bar{r}(h) = \frac{1}{2} M_2 h + \frac{2\bar{\Delta}}{h}$ (график функции $\bar{r}(h)$ для случая, рассмотренного ниже в примере 12.4, приведен на рис. 12.2).

Выберем оптимальное значение шага h , при котором величина $\bar{r}(h)$ достигает минимального значения. Приравнивая производную $\bar{r}'(h) = \frac{1}{2} M_2 - \frac{2\bar{\Delta}}{h^2}$ к нулю, получаем значение $h_{\text{опт}} = 2\sqrt{\bar{\Delta}/M_2}$, которому отвечает величина $\bar{r}_{\min} = \bar{r}(h_{\text{опт}}) = 2\sqrt{\bar{\Delta} M_2}$.

Таким образом, при использовании формулы (12.1) для вычисления производной функции f , заданной с погрешностью, следует обратить особое внимание на выбор шага h . Однако даже при оптимальном выборе шага полная погрешность окажется величиной, пропорциональной лишь $\sqrt{\bar{\Delta}}$.

Формулы для вычисления производных порядка $k > 1$ обладают еще большей чувствительностью к ошибкам задания функций. Поэтому значения производных высокого порядка, найденные с помощью таких формул, могут быть очень неточными.

Пример 12.4. Рассмотрим результаты применения формулы (12.1) с разными значениями шага h для вычисления производной функции $f(x) = e^x$ в точке $x = 1$.

В табл. 12.4 приведены значения приближений $(f'')^*$ к $f'(1) = e \approx 2.71828$, полученные на 6-разрядном десятичном компьютере и отвечающие значениям $h = 10^{-1}, 10^{-2}, \dots, 10^{-6}$. Для удобства анализа указаны также значения погрешностей r^* .

Таблица 12.4

h	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}	10^{-6}
$(f'')^*$	2.85890	2.73200	2.72000	2.70000	3.00000	0.00000
r^*	$-14 \cdot 10^{-2}$	$-14 \cdot 10^{-3}$	$-17 \cdot 10^{-4}$	$18 \cdot 10^{-3}$	$-28 \cdot 10^{-2}$	$27 \cdot 10^{-1}$

Из таблицы видно, что погрешность r^* с уменьшением h сначала убывает по модулю, а затем начинает резко возрастать. При $h = 10^{-6}$ значения функции в точках $x = 1 + h$ и $x = 1$, найденные с шестью знаками мантиссы, совпадают, и поэтому вычисление по формуле (12.1) дает приближение к $f'(1)$, равное нулю. Ясно, что при $h < 10^{-6}$ будет получаться тот же результат.

На рис. 12.2 точками помечены значения модуля погрешности r^* , отвечающие различным h из диапазона $10^{-4} \leq h \leq 10^{-2}$. Сплошной

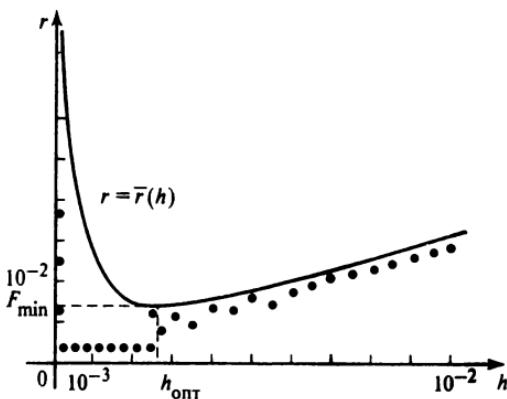


Рис. 12.2

линией изображен график функции $\bar{r}(h) = 14h + 10^{-5}/h$, являющейся верхней оценкой для $|r^*(h)|$ (в данном случае $M_2 = \max_{[1, 1.01]} |e^x| \leq 2.8$ и $\bar{\Delta} = 5 \cdot 10^{-6}$). Отметим, что хотя реальное значение погрешности r^* оказывается меньше получаемого с помощью оценки \bar{r} , все же функция $\bar{r}(h)$ правильно отражает основные особенности поведения погрешности. ▲

§ 12.4. Дополнительные замечания

1. Формулы численного дифференцирования применяются и для приближенного вычисления частных производных функций многих переменных. Для их построения используются различные приемы. Среди них последовательное применение одномерных формул численного дифференцирования и дифференцирование интерполяционных формул. Распространенным является и метод неопределенных коэффициентов [9] (иногда он используется и в одномерном случае).

2. При наличии в значениях функции случайных ошибок нередко применяют некоторые процедуры предварительного сглаживания. В последнее время получила распространение группа методов численного дифференцирования, в которых используются идеи регуляризации (см. [97]).

3. Дополнительную информацию о методах численного дифференцирования можно найти, например в [51].

ЧИСЛЕННОЕ ИНТЕГРИРОВАНИЕ

§ 13.1. Простейшие квадратурные формулы

1. Постановка задачи. В прикладных исследованиях часто возникает необходимость вычисления значения определенного интеграла

$$I = \int_a^b f(x) dx. \quad (13.1)$$

Этот интеграл может выражать площадь, объем, работу переменной силы и т.д.

Если функция $f(x)$ непрерывна на отрезке $[a, b]$ и ее первообразную $F(x)$ удается выразить через известные функции, то для вычисления интеграла (13.1) можно воспользоваться *формулой Ньютона—Лейбница*¹:

$$\int_a^b f(x) dx = F(b) - F(a). \quad (13.2)$$

К сожалению, в подавляющем большинстве случаев получить значение определенного интеграла с помощью формулы (13.2) или других аналитических методов не удается.

Пример 13.1. Интеграл $\int_0^x e^{-t^2} dt$ широко используется при исследо-

вании процессов теплообмена и диффузии, в статистической физике и теории вероятностей. Однако его значение не может быть выражено в виде конечной комбинации элементарных функций. ▲

Заметим, что даже тогда, когда удается получить первообразную функцию $F(x)$ в аналитической форме, значительные усилия, затраченные на это, часто оказываются чрезмерно высокой платой за окончательный результат. Добавим еще, что вычисления интеграла в этих случаях по формуле (13.2), как правило, приводят к громоздким (а часто — и приближенным) вычислениям. Следует отметить также, что часто

¹ Готфрид Вильгельм Лейбниц (1646—1716) — немецкий математик, физик и философ. Один из создателей дифференциального и интегрального исчислений.

найти точное значение интеграла (13.1) просто невозможно. Например, это имеет место, когда функция $f(x)$ задается таблицей своих значений.

Обычно для вычисления значения определенного интеграла применяют специальные численные методы. Наиболее широко используют на практике *квадратурные формулы* — приближенные равенства вида

$$\int_a^b f(x) dx \approx \sum_{i=0}^N A_i f(\bar{x}_i). \quad (13.3)$$

Здесь \bar{x}_i — некоторые точки из отрезка $[a, b]$ — узлы *квадратурной формулы*; A_i — числовые коэффициенты, называемые *весами квадратурной формулы*; $N \geq 0$ — целое число. Сумма $\sum_{i=0}^N A_i f(\bar{x}_i)$, которая принимается за приближенное значение интеграла, называется *квадратурной суммой*. Величина $R = \int_a^b f(x) dx - \sum_{i=0}^N A_i f(\bar{x}_i)$ называется *погрешностью* (или *остаточным членом*) *квадратурной формулы*.

Будем говорить, что квадратурная формула (13.3) *точна для многочленов степени m* , если для любого многочлена степени не выше m эта формула дает точное значение интеграла, т.е.

$$\int_a^b P_m(x) dx = \sum_{i=0}^N A_i P_m(\bar{x}_i).$$

При оценке эффективности квадратурных формул часто исходят из того, что наиболее трудоемкой операцией при вычислении по формуле (13.3) является нахождение значения функции f . Поэтому среди двух формул, позволяющих вычислить интеграл с заданной точностью ε , более эффективной считается та, в которой используется меньшее число узлов.

Выведем простейшие квадратурные формулы, используя наглядные геометрические соображения. Будем интерпретировать интеграл (13.1) как площадь криволинейной трапеции, ограниченной графиком функции $y = f(x)$ (при $f(x) \geq 0$), осью абсцисс и прямыми $x = a$, $x = b$ (рис. 13.1, а).

Разобьем отрезок $[a, b]$ на элементарные отрезки $[x_{i-1}, x_i]$ точками $a = x_0 < x_1 < \dots < x_n = b$. Интеграл I разобьется при этом на сумму элементарных интегралов:

$$I = \sum_{i=1}^n I_i, \quad (13.4)$$

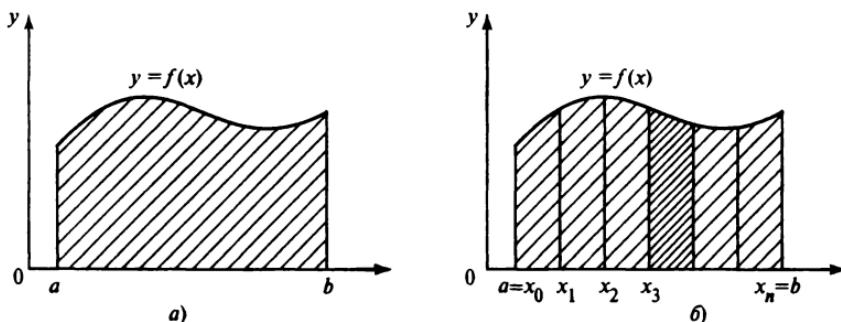


Рис. 13.1

где $I_i = \int_{x_{i-1}}^{x_i} f(x) dx$, что соответствует разбиению площади исходной криволинейной трапеции на сумму площадей элементарных криволинейных трапеций (рис. 13.1, б).

Введем обозначения: $f_i = f(x_i)$, $f_{i-1/2} = f(x_{i-1/2})$, где $x_{i-1/2} = (x_{i-1} + x_i)/2$ — середина элементарного отрезка. Для простоты шаг $h = x_i - x_{i-1}$ будем считать постоянным.

2. Формула прямоугольников. Заменим приближенно площадь элементарной криволинейной трапеции площадью прямоугольника, основанием которого является отрезок $[x_{i-1}, x_i]$, а высота равна значению $f_{i-1/2}$ (на рис. 13.2, а через $N_{i-1/2}$ обозначена точка с координа-

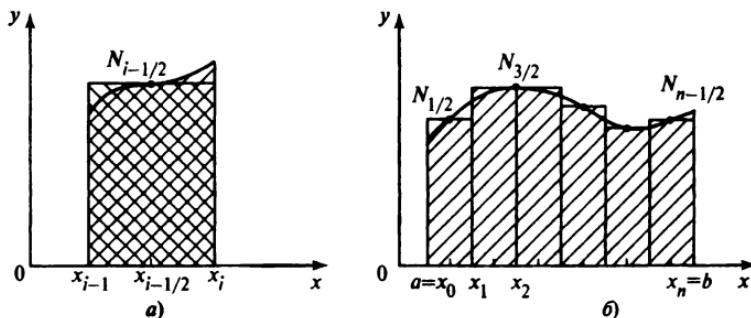


Рис. 13.2

тами $(x_{i-1/2}, f_{i-1/2})$). Так мы приходим к **элементарной квадратурной формуле прямоугольников**:

$$I_i \approx h f_{i-1/2}. \quad (13.5)$$

Производя такую замену для всех элементарных криволинейных трапеций, получаем **составную квадратурную формулу прямоугольников**:

$$I \approx I_{\text{пп}}^h = h(f_{1/2} + f_{3/2} + \dots + f_{n-1/2}) = h \sum_{i=1}^n f_{i-1/2}. \quad (13.6)$$

Эта формула соответствует приближенной замене площади исходной криволинейной трапеции площадью ступенчатой фигуры, изображенной на рис. 13.2, б.

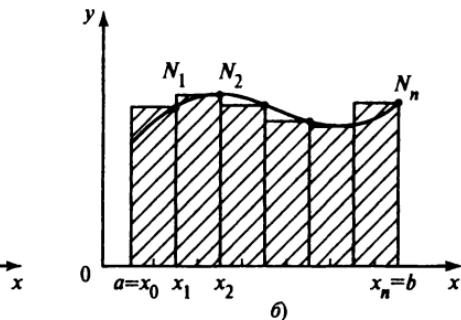
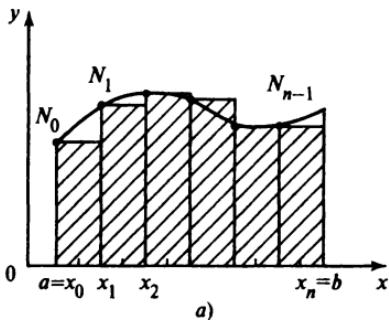


Рис. 13.3

Замечание. Иногда используют формулы

$$I \approx h \sum_{i=0}^{n-1} f_i, \quad (13.7)$$

$$I \approx h \sum_{i=1}^n f_i, \quad (13.8)$$

называемые соответственно **составными квадратурными формулами левых и правых прямоугольников**. Геометрические иллюстрации приведены на рис. 13.3, а и б. В соответствии с этим формулу (13.6) иногда называют **составной квадратурной формулой центральных прямоугольников**.

3. Формула трапеций. Соединив отрезком точки $N_{i-1}(x_{i-1}, f_{i-1})$ и $N_i(x_i, f_i)$ на графике функции $y = f(x)$, получим трапецию (рис. 13.4, а).

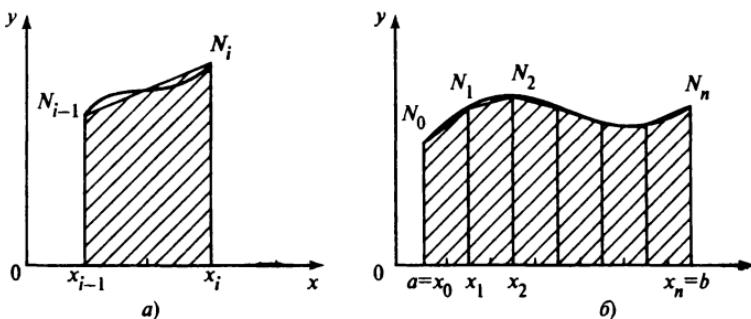


Рис. 13.4

Заменим теперь приближенно площадь элементарной криволинейной трапеции площадью построенной фигуры. Тогда получим *элементарную квадратурную формулу трапеций*:

$$I_i \approx \frac{h}{2} (f_{i-1} + f_i). \quad (13.9)$$

Пользуясь этой формулой при $i = 1, \dots, n$, выводим *составную квадратурную формулу трапеций*:

$$I \approx I_{\text{тр}}^h = h \left[\frac{f_0}{2} + f_1 + f_2 + \dots + f_{n-1} + \frac{f_n}{2} \right] = h \left[\frac{f_0 + f_n}{2} + \sum_{i=1}^{n-1} f_i \right]. \quad (13.10)$$

Эта формула соответствует приближенной замене площади исходной криволинейной трапеции площадью фигуры, ограниченной ломаной линией, проходящей через точки N_0, N_1, \dots, N_n (рис. 13.4, б).

4. Формула Симпсона¹. Если площадь элементарной криволинейной трапеции заменить площадью фигуры, расположенной под параболой, проходящей через точки $N_{i-1}, N_{i-1/2}$ и N_i (рис. 13.5, а), то получим

приближенное равенство $I_i \approx \int_{x_{i-1}}^{x_i} P_2(x) dx$. Здесь $P_2(x)$ — интерполяционный многочлен второй степени с узлами $x_{i-1}, x_{i-1/2}, x_i$. Как нетрудно убедиться², верна формула

$$P_2(x) = f_{i-1/2} + \frac{f_i - f_{i-1}}{h} (x - x_{i-1/2}) + \frac{f_i - 2f_{i-1/2} + f_{i-1}}{h^2/2} (x - x_{i-1/2})^2.$$

¹ Томас Симпсон (1710—1761) — английский математик.

² Проверку того, что $P_2(x_{i-1}) = f_{i-1}$, $P_2(x_{i-1/2}) = f_{i-1/2}$, $P_2(x_i) = f_i$, рекомендуем провести самостоятельно.

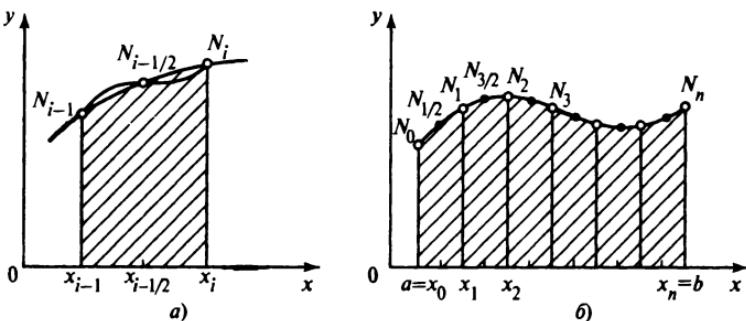


Рис. 13.5

Ее интегрирование приводит к равенству

$$\begin{aligned} \int_{x_{i-1}}^{x_i} P_2(x) dx &= hf_{i-1/2} + \frac{f_i - f_{i-1}}{h} \int_{x_{i-1}}^{x_i} (x - x_{i-1/2}) dx + \\ &+ \frac{f_i - 2f_{i-1/2} + f_{i-1}}{h^2/2} \int_{x_{i-1}}^{x_i} (x - x_{i-1/2})^2 dx = \\ &= hf_{i-1/2} + \frac{h}{6} (f_{i-1} - 2f_{i-1/2} + f_i) = \frac{h}{6} (f_{i-1} + 4f_{i-1/2} + f_i). \end{aligned}$$

Таким образом, выведена элементарная квадратурная формула Симпсона¹:

$$I_i \approx \frac{h}{6} (f_{i-1} + 4f_{i-1/2} + f_i). \quad (13.11)$$

Применяя эту формулу на каждом элементарном отрезке, выводим составную квадратурную формулу Симпсона:

$$\begin{aligned} I \approx I_C^h &= \frac{h}{6} (f_0 + 4f_{1/2} + 2f_1 + 4f_{3/2} + 2f_2 + \dots + 2f_{n-1} + 4f_{n-1/2} + f_n) = \\ &= \frac{h}{6} \left(f_0 + f_n + 4 \sum_{i=1}^n f_{i-1/2} + 2 \sum_{i=1}^{n-1} f_i \right) \quad (13.12) \end{aligned}$$

Замечание 1. Учитывая геометрическую интерпретацию формулы Симпсона, ее иногда называют *формулой парабол*.

¹ Формула Симпсона впервые была получена Кавальери в 1639 г. Затем в 1668 г. она была переоткрыта Джеймсом Грегори. Томас Симпсон вывел ее в 1743 г.

Замечание 2. Когда число элементарных отрезков разбиения четно ($n = 2m$), в формуле Симпсона можно использовать лишь узлы с целыми индексами:

$$I \approx \frac{h}{3} \left(f_0 + f_{2m} + 4 \sum_{i=1}^m f_{2i-1} + 2 \sum_{i=1}^{m-1} f_{2i} \right).$$

При выводе этой формулы роль элементарного отрезка играет отрезок $[x_{2i-2}, x_{2i}]$ длины $2h$.

5. Оценка погрешности. Оценим погрешность выведенных квадратурных формул в предположении, что подынтегральная функция f достаточно гладкая. Как и в предыдущих главах, будем использовать обозначение $M_k = \max_{[a, b]} |f^{(k)}(x)|$.

Теорема 13.1. Пусть функция f дважды непрерывно дифференцируема на отрезке $[a, b]$. Тогда для составных квадратурных формул прямоугольников и трапеций справедливы следующие оценки погрешности:

$$|I - I_{\text{пп}}^h| \leq \frac{M_2(b-a)}{24} h^2, \quad (13.13)$$

$$|I - I_{\text{тр}}^h| \leq \frac{M_2(b-a)}{12} h^2. \quad (13.14)$$

Доказательство. Выведем сначала оценку (13.13). Представим погрешность $R = I - I_{\text{пп}}^h$ формулы прямоугольников в виде

$$R = \int_a^b f(x) dx - h \sum_{i=1}^n f_{i-1/2} = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} (f(x) - f(x_{i-1/2})) dx.$$

Используя формулу Тейлора

$$f(x) = f(x_{i-1/2}) + f'(x_{i-1/2})(x - x_{i-1/2}) + \frac{f''(\xi)}{2} (x - x_{i-1/2})^2,$$

где $x \in [x_{i-1}, x_i]$, $\xi = \xi(x) \in [x_{i-1}, x_i]$, имеем

$$R_i = \int_{x_{i-1}}^{x_i} (f(x) - f(x_{i-1/2})) dx = \frac{1}{2} \int_{x_{i-1}}^{x_i} f''(\xi(x))(x - x_{i-1/2})^2 dx,$$

$$|R_i| \leq \frac{M_2}{2} \int_{x_{i-1}}^{x_i} (x - x_{i-1/2})^2 dx = \frac{M_2}{6} (x - x_{i-1/2})^3 \Big|_{x_{i-1}}^{x_i} = \frac{M_2}{24} h^3.$$

Так как $R = \sum_{i=1}^n R_i$, то $|R| \leq \sum_{i=1}^n \frac{M_2}{24} h^3 = \frac{M_2}{24} h^3 n$. Замечая, что $nh = b - a$, приходим к оценке (13.13).

Для вывода оценки (13.14) воспользуемся тем, что отрезок, соединяющий точки N_{i-1} и N_i представляет собой график интерполяционного многочлена первой степени $P_1(x) = f_{i-1} \frac{x_i - x}{h} + f_i \frac{x - x_{i-1}}{h}$.

Поэтому для элементарной формулы трапеций верно равенство

$$R_i = \int_{x_{i-1}}^{x_i} f(x) dx - \frac{h}{2} (f_{i-1} + f_i) = \int_{x_{i-1}}^{x_i} (f(x) - P_1(x)) dx.$$

Используя оценку (11.28) погрешности линейной интерполяции, имеем

$$|R_i| \leq \int_{x_{i-1}}^{x_i} \frac{M_2}{2} (x - x_{i-1})(x_i - x) dx = \frac{M_2}{12} h^3.$$

Следовательно, для $R = I - I_{\text{тр}}^h$ справедлива оценка

$$|R| \leq \sum_{i=1}^n |R_i| \leq \frac{M_2}{12} h^3 n = \frac{M_2(b-a)}{12} h^2. \blacksquare$$

Приведем теперь без доказательства теорему об оценке погрешности формулы Симпсона.

Теорема 13.2. Пусть функция f имеет на отрезке $[a, b]$ непрерывную производную четвертого порядка $f^{(4)}$. Тогда для формулы Симпсона (13.12) справедлива оценка погрешности

$$|I - I_C^h| \leq \frac{M_4(b-a)}{2880} h^4. \blacksquare \quad (13.15)$$

Замечание 1. Оценки (13.13)–(13.15) означают, что формулы прямоугольников и трапеций имеют второй порядок точности относительно h , а формула Симпсона — четвертый порядок точности. Из тех же оценок следует, что формулы прямоугольников и трапеций точны для многочленов первой степени, а формула Симпсона — для многочленов третьей степени.

Замечание 2. Формулы (13.7) и (13.8) имеют лишь первый порядок точности (абсолютная погрешность каждой из формул не превышает $0.5M_1(b-a)h$) и поэтому для вычисления интегралов на практике они используются крайне редко.

Пример 13.2. Вычислим значение интеграла $\int_0^1 e^{-x^2} dx$, используя

квадратурные формулы прямоугольников, трапеций и Симпсона с шагом $h = 0.1$.

Сначала составим таблицу значений функции $y = e^{-x^2}$ (табл. 13.1).

Таблица 13.1

x	e^{-x^2}	x	e^{-x^2}
0.00	1.0000000	0.55	0.7389685
0.05	0.9975031	0.60	0.6976763
0.10	0.9900498	0.65	0.6554063
0.15	0.9777512	0.70	0.6126264
0.20	0.9607894	0.75	0.5697828
0.25	0.9394131	0.80	0.5272924
0.30	0.9139312	0.85	0.4855369
0.35	0.8847059	0.90	0.4448581
0.40	0.8521438	0.95	0.4055545
0.45	0.8166865	1.00	0.3678794
0.50	0.7788008		

Производя вычисления по формулам (13.6), (13.10), (13.12), получаем $I_{\text{пп}}^h = 0.74713088$, $I_{\text{tp}}^h = 0.74621079$, $I_C^h = 0.74682418$.

Оценим погрешность каждого из полученных значений, используя неравенства (13.13)–(13.15). Вычислим $f^{(2)}(x) = (4x^2 - 2)e^{-x^2}$. Как нетрудно видеть $|f^{(2)}| \leq M_2 = 2$. Следовательно,

$$|I - I_{\text{пп}}^h| \leq \frac{2 \cdot 1}{24} (0.1)^2 \approx 0.84 \cdot 10^{-3}, \quad |I - I_{\text{tp}}^h| \leq \frac{2 \cdot 1}{12} (0.1)^2 \approx 1.7 \cdot 10^{-3}.$$

Далее $f^{(4)}(x) = (16x^4 - 48x^2 + 12)e^{-x^2}$, $|f^{(4)}| \leq 12$. Поэтому

$$|I - I_C^h| \leq \frac{12 \cdot 1}{2880} (0.1)^4 \approx 0.42 \cdot 10^{-6}.$$

Таким образом, из вычислений по формуле прямоугольников с учетом погрешности следует, что $I = 0.747 \pm 0.001$; по формуле трапеций — что $I = 0.746 \pm 0.002$; по формуле Симпсона — что $I = 0.7468242 \pm 0.0000005$. ▲

6. Случай переменного шага. Приведем составные квадратурные формулы прямоугольников, трапеций и Симпсона с переменным шагом $h_i = x_i - x_{i-1}$ и $h_{i+1/2} = (h_{i+1} + h_i)/2$.

$$I \approx I_{\text{np}}^h = \sum_{i=1}^n f_{i-1/2} h_i,$$

$$I \approx I_{\text{tp}}^h = \sum_{i=1}^n \frac{f_{i-1} + f_i}{2} h_i = f_0 \frac{h_1}{2} + \sum_{i=1}^{n-1} f_i h_{i+1/2} + f_n \frac{h_n}{2},$$

$$I \approx I_{\text{C}}^h = \sum_{i=1}^n \frac{f_{i-1} + 4f_{i-1/2} + f_i}{6} h_i.$$

Вывод этих формул и их геометрический смысл остаются теми же, что и для постоянного шага. Теоремы об оценках погрешности также останутся справедливыми, если в неравенствах (13.13)–(13.15) заменить h на $h_{\max} = \max_{1 \leq i \leq n} h_i$.

§ 13.2. Квадратурные формулы интерполяционного типа

Для приближенного вычисления определенных интегралов часто используется следующий естественный для методов приближения функций прием. Подынтегральную функцию f аппроксимируют на отрезке $[a, b]$ некоторой функцией g , интеграл от которой легко вычисляется, а затем полагают

$$\int_a^b f(x) dx \approx \int_a^b g(x) dx. \quad (13.16)$$

Точность формулы (13.16) можно повышать за счет усложнения метода глобальной аппроксимации. Однако чаще используется другой подход. Интеграл I представляют в виде суммы (13.4) интегралов по элементарным отрезкам $[x_{i-1}, x_i]$. На каждом таком i -м отрезке функцию f аппроксимируют некоторой легко интегрируемой функцией g_i . В результате получается составная формула

$$\int_a^b f(x) dx \approx \sum_{i=1}^n \int_{x_{i-1}}^{x_i} g_i(x) dx.$$

1. Вывод квадратурных формул интерполяционного типа. Рассмотрим более подробно этот подход в случае, когда аппроксимация осуществляется с помощью интерполяционного многочлена. Зафиксируем некоторые значения $t_0, t_1, \dots, t_m \in [-1, 1]$. Аппроксимируем функцию $f(x)$ на i -м элементарном отрезке $[x_{i-1}, x_i]$ интерполяционным многочленом $P_{m,i}(x)$ с узлами интерполяции $z_j^{(i)} = x_{i-1/2} + t_j h_i/2$, $j = 0, 1, 2, \dots, m$. В случае, когда все значения t_j различны, можно воспользоваться записью интерполяционного многочлена в форме Лагранжа:

$$P_{m,i}(x) = \sum_{j=0}^m f(z_j^{(i)}) l_{m,j}^{(i)}(x), \quad l_{m,j}^{(i)}(x) = \prod_{\substack{k=0 \\ k \neq j}}^m \frac{x - z_k^{(i)}}{z_j^{(i)} - z_k^{(i)}}. \quad (13.17)$$

Используя замену переменной $x = x_{i-1/2} + th_i/2$, вычислим интеграл от $P_{m,i}$ на отрезке $[x_{i-1}, x_i]$:

$$\begin{aligned} \int_{x_{i-1}}^{x_i} P_{m,i}(x) dx &= \sum_{j=0}^m f(z_j^{(i)}) \int_{x_{i-1}}^{x_i} l_{m,j}^{(i)}(x) dx = \\ &= h_i \sum_{j=0}^m a_j f(x_{i-1/2} + t_j h_i/2), \quad a_j = \frac{1}{2} \int_{-1}^1 \prod_{\substack{k=0 \\ k \neq j}}^m \frac{(t - t_k)}{(t_j - t_k)} dt. \end{aligned}$$

Приближенная замена интеграла I суммой $I^h = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} P_{m,i}(x) dx$

приводит к следующей *составной квадратурной формуле интерполяционного типа*:

$$I \approx I^h = \sum_{i=1}^n h_i \sum_{j=0}^m a_j f(x_{i-1/2} + t_j h_i/2). \quad (13.18)$$

Замечание. Квадратурные формулы интерполяционного типа, построенные на основе равноотстоящих значений t_0, t_1, \dots, t_m называют *формулами Ньютона—Котеса*¹.

Рассмотренные в § 13.1 простейшие квадратурные формулы являются формулами интерполяционного типа; более того, они относятся к классу формул Ньютона—Котеса.

¹ Роджер Котес (1682—1716) — английский математик, друг и ученик И. Ньютона.

Приведем квадратурные формулы Ньютона—Котеса, отвечающие использованию многочленов степени $m = 1, 2, 3, 4, 5, 6$ соответственно.

$$I \approx \sum_{i=1}^n \frac{h_i}{2} (f(x_{i-1}) + f(x_i)), \quad m = 1 — \text{формула трапеций};$$

$$I \approx \sum_{i=1}^n \frac{h_i}{6} (f(x_{i-1}) + 4f(x_{i-1/2}) + f(x_i)), \quad m = 2 — \text{формула Симпсона};$$

$$I \approx \sum_{i=1}^n \frac{h_i}{8} \left(f(x_{i-1}) + 3f\left(x_{i-1} + \frac{h_i}{3}\right) + 3f\left(x_i - \frac{h_i}{3}\right) + f(x_i) \right),$$

$m = 3 — \text{правило } 3/8;$

$$I \approx \sum_{i=1}^n \frac{h_i}{90} \left(7f(x_{i-1}) + 32f\left(x_{i-1} + \frac{h_i}{4}\right) + 12f(x_{i-1/2}) + 32f\left(x_i - \frac{h_i}{4}\right) + 7f(x_i) \right),$$

$m = 4 — \text{формула Милна (формула Боде);}$

$$I \approx \sum_{i=1}^n \frac{h_i}{288} \left(19f(x_{i-1}) + 75f\left(x_{i-1} + \frac{h_i}{5}\right) + 50\left(x_{i-1} + \frac{2h_i}{5}\right) + 50f\left(x_i - \frac{2h_i}{5}\right) + 75f\left(x_i - \frac{h_i}{5}\right) + 19f(x_i) \right),$$

$m = 5;$

$$I \approx \sum_{i=1}^n \frac{h_i}{840} \left(41f(x_{i-1}) + 216f\left(x_{i-1} + \frac{h_i}{6}\right) + 27f\left(x_{i-1} + \frac{h_i}{3}\right) + 272f(x_{i-1/2}) + 27f\left(x_i - \frac{h_i}{3}\right) + 216f\left(x_i - \frac{h_i}{6}\right) + 41f(x_i) \right),$$

$m = 6 — \text{формула Вэддла.}$

2. Оценка погрешности. Приведем теорему об оценке погрешности формулы (13.18).

Теорема 13.3. Пусть функция f имеет на отрезке $[a, b]$ непрерывную производную порядка $m+1$. Тогда для погрешности квадратурной формулы (13.18) справедлива оценка

$$|I - I^h| \leq C_m M_{m+1}(b-a) h_{\max}^{m+1}, \quad (13.19)$$

где

$$C_m = \frac{1}{2^{m+2}(m+1)!} \int_{-1}^1 |\bar{\omega}_{m+1}(t)| dt, \quad \bar{\omega}_{m+1}(t) = \prod_{k=0}^m (t - t_k).$$

Доказательство. Представим погрешность $R = I - I^h$ формулы (13.18) в виде

$$R = \int_a^b f(x) dx - \sum_{i=1}^n \int_{x_{i-1}}^{x_i} P_{m,i}(x) dx = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} (f(x) - P_{m,i}(x)) dx.$$

Пользуясь оценкой (11.26) погрешности интерполяции, в данном случае принимающей вид

$$|f(x) - P_{m,i}(x)| \leq \frac{M_{m+1}}{(m+1)!} \left| \prod_{k=0}^m (x - z_k^{(i)}) \right|,$$

и производя замену переменной $x = x_{i-1/2} + th_i/2$, получаем цепочку неравенств

$$\begin{aligned} |R| &\leq \sum_{i=1}^n \int_{x_{i-1}}^{x_i} |f(x) - P_{m,i}(x)| dx \leq \sum_{i=1}^n \frac{M_{m+1}}{(m+1)!} \int_{x_{i-1}}^{x_i} \left| \prod_{k=0}^m (x - z_k^{(i)}) \right| dx = \\ &= \sum_{i=1}^n \left(\frac{h_i}{2} \right)^{m+2} \frac{M_{m+1}}{(m+1)!} \int_{-1}^1 |\bar{\omega}_{m+1}(t)| dt = C_m M_{m+1} \sum_{i=1}^n h_i^{m+2}. \end{aligned}$$

Учитывая, что $\sum_{i=1}^n h_i^{m+2} \leq h_{\max}^{m+1} \sum_{i=1}^n h_i = h_{\max}^{m+1} (b-a)$, приходим к оценке (13.19). ■

Замечание 1. Теорема остается справедливой и в случае, когда для построения квадратурной формулы используется интерполяция с кратными узлами (т.е. когда некоторые из значений t_0, t_1, \dots, t_m совпадают).

Замечание 2. Как следует из оценки (13.19), квадратурные формулы интерполяционного типа (13.18) точны для многочленов степени m .

Отметим, что для некоторых симметричных квадратурных формул оценка погрешности (13.19) является грубой и не отражает истинный порядок их точности. Например, согласно этой оценке формулы прямоугольников и Симпсона должны иметь лишь первый и третий порядок точности соответственно, а в действительности они имеют на единицу больший порядок точности (см. теоремы 13.1, 13.2 и замечание 1 на с. 426). Поясним причину этого явления на примере формулы прямоугольников.

Аппроксимируем функцию f на отрезке $[x_{i-1}, x_i]$ интерполяционным многочленом $P_{1,i}(x) = f(x_{i-1/2}) + f'(x_{i-1/2})(x - x_{i-1/2})$ с кратным узлом $x_{i-1/2}$. Интегрирование многочлена $P_{1,i}(x)$ дает значение $I_i \approx$

$$\approx \int_{x_{i-1}}^{x_i} P_{1,i}(x) dx = h_i f_{i-1/2}, \text{ совпадающее со значением, получаемым по}$$

элементарной формуле прямоугольников (13.5). Геометрическая иллюстрация, приведенная на рис. 13.6, показывает, что формулу (центральных) прямоугольников с равным основанием можно было бы назвать и формулой трапеций. Действительно, площадь $f_{i-1/2} \cdot h_i$ элементарного прямоугольника $ABCD$ совпадает с площадью трапеции $AEFD$, одна из сторон EF которой является касательной к кривой $y = f(x)$, проведенной в точке $N_{i-1/2}$ (ведь площади треугольников $EBN_{i-1/2}$ и $N_{i-1/2}FC$ равны).

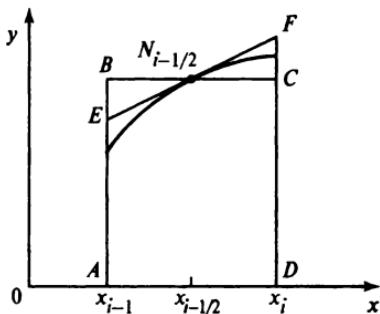


Рис. 13.6

Таким образом, можно считать, что формула прямоугольников построена на основе интерполяции функции f многочленом первой степени с кратным узлом. Теперь из теоремы 13.3 с учетом замечания 1 следует, что эта формула действительно имеет второй порядок точности.

Формула Симпсона может быть получена (аналогично) на основе интерполяции многочленом третьей степени с узлами $x_{i-1}, x_{i-1/2}, x_i$, центральный из которых является кратным. Следовательно, эту формулу с достаточным основанием можно было бы называть и формулой кубических парабол. Закономерно, что формула Симпсона имеет четвертый порядок точности.

Приведем оценки погрешности формул Ньютона—Котеса.

$$|I - I^h| \leq \frac{M_2(b-a)}{12} h_{\max}^2, m = 1 — \text{формула трапеций};$$

$$|I - I^h| \leq \frac{M_4(b-a)}{2880} h_{\max}^4, \quad m = 2 \text{ — формула Симпсона;}$$

$$|I - I^h| \leq \frac{M_4(b-a)}{6480} h_{\max}^4, \quad m = 3 \text{ — правило } 3/8;$$

$$|I - I^h| \leq \frac{M_6(b-a)}{1935360} h_{\max}^6, \quad m = 4 \text{ — формула Милна (формула Боде);}$$

$$|I - I^h| \leq \frac{11M_6(b-a)}{37800000} h_{\max}^6, \quad m = 5;$$

$$|I - I^h| \leq \frac{M_8(b-a)}{1567641600} h_{\max}^8, \quad m = 6 \text{ — формула Вэддла.}$$

3. Обусловленность квадратурных формул интерполяционного типа. При вычислении интегралов, как правило, приходится использовать не точные значения $f(x)$ подынтегральной функции, а приближенные значения $f^*(x)$. Напомним (см. § 3.1), что задача вычисления определенного интеграла от приближенно заданной функции является устойчивой. В предположении, что $|f(x) - f^*(x)| \leq \bar{\Delta}(f^*)$ для всех $x \in [a, b]$, справедлива оценка

$$\Delta(I^*) = |I - I^*| = \left| \int_a^b f(x) dx - \int_a^b f^*(x) dx \right| \leq (b-a)\bar{\Delta}(f^*),$$

указывающая на то, что абсолютное число обусловленности этой задачи равно $b-a$, т.е. длине отрезка интегрирования.

Какова же чувствительность квадратурной формулы (13.3) к погрешностям задания функции f ? Заметим, что

$$\left| \sum_{i=0}^N A_i f(\bar{x}_i) - \sum_{i=0}^N A_i f^*(\bar{x}_i) \right| \leq \left[\sum_{i=0}^N |A_i| \right] \bar{\Delta}(f^*).$$

Таким образом, квадратурная формула устойчива к погрешностям задания функции и ее число обусловленности v_Δ равно $\sum_{i=0}^N |A_i|$.

Заметим, что все квадратурные формулы интерполяционного типа точны для многочленов нулевой степени и поэтому $b-a = \int_a^b 1 \cdot dx = \sum_{i=0}^N A_i$. Следовательно, если все веса A_i квадратурной формулы

интерполяционного типа положительны, то ее число обусловленности v_Δ совпадает с $b - a$. Чувствительность такой формулы к погрешностям адекватна чувствительности вычисляемого интеграла. Если же среди весов A_i имеются отрицательные, то $v_\Delta = \sum_{i=0}^N |A_i| > \sum_{i=0}^N A_i = b - a$.

Известно, что при больших значениях m среди весов квадратурной формулы (13.18) появляются отрицательные и значение числа обусловленности v_Δ становится большим. Например, для формул Ньютона—Котеса $v_\Delta \approx 3.1(b - a)$ при $m = 10$, $v_\Delta \approx 8.3(b - a)$ при $m = 20$, $v_\Delta \approx 560(b - a)$ при $m = 30$. В силу плохой обусловленности эти формулы уже при $m \geq 10$ используются весьма редко.

§ 13.3. Квадратурные формулы Гаусса

1. Построение квадратурных формул Гаусса. Из результатов

§ 13.2 следует, что квадратурная формула $\int\limits_a^b f(x) dx \approx \sum_{i=0}^N A_i f(x_i)$,

построенная интегрированием интерполяционного многочлена степени N с фиксированными узлами x_0, x_1, \dots, x_N , точна для всех многочленов степени N . Однако, если имеется свобода в выборе узлов, то можно распорядиться ею так, чтобы получить формулу, точную для всех многочленов некоторой степени, превышающей N .

Поставим следующую задачу: при заданном числе $N + 1$ узлов построить квадратурную формулу, точную для многочленов наиболее высокой степени. Формулы, удовлетворяющие этому условию, принято называть *квадратурными формулами Гаусса*. Как правило, сначала строят формулы Гаусса

$$\int\limits_{-1}^1 f(t) dt \approx G_{N+1} = \sum_{i=0}^N a_i f(t_i) \quad (13.20)$$

для стандартного отрезка $[-1, 1]$. Затем с помощью замены переменной $x = \frac{a+b}{2} + \frac{b-a}{2} t$ осуществляют переход к формулам интегрирования на произвольном отрезке:

$$\int\limits_a^b f(x) dx \approx \frac{b-a}{2} \sum_{i=0}^N a_i f\left[\frac{a+b}{2} + \frac{b-a}{2} t_i\right]. \quad (13.21)$$

Заметим, что формула (13.20) точна для многочленов степени m тогда и только тогда, когда она точна для функций $f(t) = 1, t, t^2, \dots, t^m$.

Это эквивалентно тому, что узлы t_i и весы a_i формулы (13.20) должны удовлетворять системе нелинейных уравнений

$$\sum_{i=0}^N a_i t_i^k = \int_{-1}^1 t^k dt = \frac{1 - (-1)^{k+1}}{k+1}, \quad k = 0, 1, \dots, m. \quad (13.22)$$

Можно показать, что система (13.22) имеет единственное решение $a_0, a_1, \dots, a_N, t_0, t_1, \dots, t_N$ (причем, $t_i \in [-1, 1]$) тогда и только тогда, когда число уравнений системы совпадает с числом неизвестных, т.е. при $m = 2N + 1$.

Пример 13.3. Построим квадратурную формулу Гаусса (13.20) с двумя узлами. В этом случае, т.е. при $N = 1, m = 3$, система (13.22) примет вид

$$\begin{aligned} a_0 + a_1 &= \int_{-1}^1 1 dt = 2, & a_0 t_0 + a_1 t_1 &= \int_{-1}^1 t dt = 0, \\ a_0 t_0^2 + a_1 t_1^2 &= \int_{-1}^1 t^2 dt = \frac{2}{3}, & a_0 t_0^3 + a_1 t_1^3 &= \int_{-1}^1 t^3 dt = 0. \end{aligned}$$

Решая ее, находим значения $a_0 = a_1 = 1, t_0 = -1/\sqrt{3}, t_1 = 1/\sqrt{3}$. Таким образом, получаем квадратурную формулу Гаусса

$$\int_{-1}^1 f(t) dt \approx f(-1/\sqrt{3}) + f(1/\sqrt{3}),$$

точную для многочленов третьей степени. ▲

Для квадратурной формулы Гаусса (13.20) справедлива следующая оценка погрешности:

$$|R| \leq \alpha_N M_{2N+2} (b-a)^{2N+3}.$$

Входящий в нее коэффициент $\alpha_N = \frac{[(N+1)!]^4}{(2N+3)[(2N+2)!]^3}$ очень быстро убывает с ростом N . Приведем, например несколько первых его значений: $\alpha_0 \approx 4 \cdot 10^{-2}, \alpha_1 \approx 2 \cdot 10^{-4}, \alpha_2 \approx 5 \cdot 10^{-7}, \alpha_3 \approx 6 \cdot 10^{-10}, \alpha_4 \approx 4 \cdot 10^{-13}$.

Можно было бы разбить отрезок интегрирования на частичные отрезки и, исходя из формулы Гаусса, построить составную формулу, имеющую порядок точности, равный $2N + 2$. Однако при интегрировании достаточно гладких функций в этом нет необходимости, так как уже при небольшом числе узлов ($4 \leq N \leq 10$) формула Гаусса обеспечивает очень высокую точность. На практике используются и формулы с десятками и сотнями узлов.

2. Узлы и веса квадратурной формулы Гаусса. Приведем значения узлов и весов квадратурной формулы Гаусса с числом узлов от 1 до 6 (табл. 13.2).

Таблица 13.2

Узлы и веса	Число узлов		
	1	2	3
t_0	0.000000000000	-0.5773502692	-0.7745966692
a_0	2.000000000000	1.0000000000	0.5555555556
t_1	—	0.5773502692	0.0000000000
a_1	—	1.0000000000	0.8888888888
t_2	—	—	0.7745966692
a_2	—	—	0.5555555556

Продолжение табл. 13.2

Узлы и веса	Число узлов		
	4	5	6
t_0	-0.8611363115	-0.9061798459	-0.9324695142
a_0	0.3478548451	0.2369268851	0.1713244924
t_1	-0.3399810436	-0.5384693101	-0.6612093864
a_1	0.6521451549	0.4786286705	0.3607615730
t_2	0.3399810436	0.0000000000	-0.2386191861
a_2	0.6521451549	0.5688888888	0.4679139346
t_3	0.8611363115	0.5384693101	0.2386191861
a_3	0.3478548451	0.4786286705	0.4679139346
t_4	—	0.9061798459	0.6612093864
a_4	—	0.2369268851	0.3607615730
t_5	—	—	0.9324695142
a_5	—	—	0.1713244924

Численные значения узлов и весов квадратурных формул Гаусса для значений N до 95 можно найти в справочнике [92].

Пример 13.4. Найдем значение интеграла $\int_0^1 e^{-x^2} dx$, используя

квадратурную формулу Гаусса с двумя, тремя и четырьмя узлами. В данном случае $a = 0$, $b = 1$, $f(x) = e^{-x^2}$ и формула (13.21) принимает вид

$$I = \int_0^1 e^{-x^2} dx \approx I_N = \frac{1}{2} \sum_{i=0}^N a_i e^{-(0.5 + 0.5t_i)^2}.$$

Взяв из табл. 13.2 значения узлов t_i и весов a_i , при $N = 1, 2, 3$, получим следующие приближения:

$$I \approx I_1 = \frac{1}{2} (e^{-(0.211324865)^2} + e^{-(0.788675135)^2}) \approx 0.7465946885;$$

$$I \approx I_2 = \frac{1}{2} (0.5555555556 \cdot e^{-(0.1127011665)^2} + 0.8888888888 \cdot e^{-0.25} + \\ + 0.5555555556 \cdot e^{-(0.887298335)^2}) \approx 0.7468145842;$$

$$I \approx I_3 = \frac{1}{2} (0.3478548451 \cdot e^{-(0.069431844)^2} + 0.6521451549 \cdot e^{-(0.330009478)^2} + \\ + 0.6521451549 \cdot e^{-(0.669990522)^2} + 0.3478548451 \cdot e^{-(0.930568156)^2}) \approx \\ \approx 0.7468244681.$$

Эти значения содержат такие абсолютные погрешности: $\Delta(I_1) \approx \approx 2 \cdot 10^{-4}$, $\Delta(I_2) \approx 10^{-5}$, $\Delta(I_3) \approx 3 \cdot 10^{-7}$. Для сравнения укажем, что значение того же интеграла, полученное в примере 13.2 по формуле прямоугольников с 10 узлами, имеет абсолютную погрешность примерно $3 \cdot 10^{-4}$, а по формуле Симпсона с 21 узлом — абсолютную погрешность примерно $5 \cdot 10^{-8}$. ▲

Замечание. Хорошо известно, что узлы квадратурных формул Гаусса G_N совпадают с корнями многочленов Лежандра $P_N(x)$, поэтому эти квадратурные формулы иногда называют формулами Гаусса—Лежандра. О связи формул Гаусса с ортогональными многочленами см., например в [9].

3. Обусловленность квадратурных формул Гаусса. Квадратурные формулы Гаусса обладают еще одним замечательным свойством: их весовые коэффициенты всегда положительны. Это свойство (как следует из рассуждений § 13.3) гарантирует хорошую обусловленность квадратурной формулы. Более того, число обусловленности равно $b - a$ и не зависит от числа узлов. Это позволяет применять на практике квадратурные формулы Гаусса с числом узлов, достигающим сотен.

§ 13.4. Апостериорные оценки погрешности. Понятие об адаптивных процедурах численного интегрирования

Применение неравенств типа (13.13)–(13.15), (13.19) для априорной оценки погрешности квадратурных формул в большинстве случаев оказывается неэффективным или вообще невозможным. Это связано как с трудностями оценивания производных подынтегральной функции f ,

так и с тем, что получаемые оценки, как правило, бывают сильно завышенными. На практике обычно используются иные подходы к оценке погрешности, позволяющие строить процедуры численного интегрирования с автоматическим выбором шага.

1. Главный член погрешности. Пусть I^h — приближенное значение интеграла $\int_a^b f(x) dx$, вычисленное по некоторой квадратурной формуле и использующее разбиение отрезка $[a, b]$ на элементарные отрезки длины h . Предположим, что для погрешности этой формулы справедливо представление¹

$$I - I^h = Ch^k + o(h^k), \quad (13.23)$$

где $C \neq 0$ и $k > 0$ — величины, не зависящие от h . Тогда величина Ch^k называется *главным членом погрешности* квадратурной формулы.

Заметим, что из неравенства (13.23) следует справедливость оценки $|I - I^h| \leq \bar{C}h^k$ с некоторой постоянной $\bar{C} > |C|$; поэтому число k представляет собой не что иное как порядок точности соответствующей квадратурной формулы.

Если подынтегральная функция f достаточно гладкая, то для каждой из составных квадратурных формул

$$I \approx I^h = \sum_{i=0}^n h \sum_{j=0}^m a_j f(x_{i-1/2} + t_j h/2) \quad (13.24)$$

существует главный член погрешности. Приведем без доказательства соответствующий результат.

Теорема 13.4. Пусть $\sum_{j=0}^m a_j = 1$ и k — минимальное среди натуральных чисел, для которых величина

$$\sigma_k = \frac{1}{2} \int_{-1}^1 t^k dt - \sum_{j=0}^m a_j t_j^k$$

¹ Напомним, что запись $\varphi(h) = o(h^k)$ (читается « $\varphi(h)$ есть o малое от h^k ») означает, что $\varphi(h)/h^k \rightarrow 0$ при $h \rightarrow 0$.

отлична от нуля. Если функция f непрерывно дифференцируема k раз на отрезке $[a, b]$, то для погрешности квадратурной формулы (13.24) справедливо представление (13.23), в котором

$$C = \frac{\sigma_k}{2^k k!} \int_a^b f^{(k)}(x) dx = \frac{\sigma_k}{2^k k!} (f^{(k-1)}(b) - f^{(k-1)}(a)). \blacksquare$$

Следствие 1. Если функция f дважды непрерывно дифференцируема на отрезке $[a, b]$, то для погрешностей составных квадратурных формул прямоугольников и трапеций справедливы следующие представления:

$$I - I_{\text{пп}}^h = C_{\text{пп}} h^2 + o(h^2), \quad C_{\text{пп}} = \frac{1}{24} \int_a^b f''(x) dx; \quad (13.25)$$

$$I - I_{\text{тр}}^h = C_{\text{тр}} h^2 + o(h^2), \quad C_{\text{тр}} = -\frac{1}{12} \int_a^b f''(x) dx. \quad (13.26)$$

Следствие 2. Если функция f четырежды непрерывно дифференцируема на отрезке $[a, b]$, то для погрешности составной квадратурной формулы Симпсона справедливо представление

$$I - I_C = C_C h^4 + o(h^4), \quad C_C = -\frac{1}{2880} \int_a^b f^{(4)}(x) dx. \quad (13.27)$$

В силу предположения (13.23) для погрешности квадратурной формулы при достаточно малом h справедливо приближенное равенство

$$I - I^h \approx Ch^k. \quad (13.28)$$

Несмотря на элементарный характер формулы (13.28), она позволяет сделать ряд важных выводов. Первый из них состоит в том, что уменьшение шага h в M раз приводит к уменьшению погрешности квадратурной формулы примерно в M^k раз. Действительно, при $h_1 = h/M$ имеем

$$I - I^{h_1} \approx Ch_1^k = \frac{1}{M^k} Ch^k \approx \frac{1}{M^k} (I - I^h).$$

В частности, уменьшение шага h в 2 раза приводит к уменьшению погрешности примерно в 2^k раз:

$$I - I^{h/2} \approx \frac{1}{2^k} Ch^k \approx \frac{1}{2^k} (I - I^h). \quad (13.29)$$

2. Правило Рунге практической оценки погрешности. Как следует из теоремы 13.4, главный член погрешности квадратурной формулы интерполяционного типа имеет вид

$$\frac{\sigma_k}{2^k k!} \int_a^b f^{(k)}(x) dx \cdot h^k = \frac{\sigma_k}{2^k k!} (f^{(k-1)}(b) - f^{(k-1)}(a)) h^k.$$

Непосредственное использование этой формулы для оценки погрешности $I - I^h$ неудобно, так как требует вычисления производных функции f . В более сложных ситуациях выражение для главного члена погрешности может оказаться существенно более громоздким. Поэтому в вычислительной практике часто применяются методы оценки погрешности, не использующие явное выражение для главного члена.

Вычитая из равенства (13.28) равенство (13.29), получаем

$$I^{h/2} - I^h \approx \frac{1}{2^k} Ch^k (2^k - 1).$$

Учитывая приближенное равенство (13.29), приходим к следующей приближенной формуле:

$$I - I^{h/2} \approx \frac{I^{h/2} - I^h}{2^k - 1}. \quad (13.30)$$

Использование этой формулы для апостериорной оценки погрешности значения $I^{h/2}$ принято называть *правилом Рунге* (или *правилом двойного пересчета*).

Замечание 1. Так как $I - I^h \approx 2^k(I - I^{h/2})$, то из (13.30) следует формула $I - I^h \approx \frac{2^k(I^{h/2} - I^h)}{2^k - 1}$, которую можно было бы использовать для приближенной оценки погрешности значения I^h . Как правило, этого не делают, поскольку среди двух вычисляемых значений интеграла I^h и $I^{h/2}$ второе является более точным и имеет смысл оценивать именно его погрешность.

Замечание 2. Заменой h на $2h$ формула (13.30) приводится к следующему виду:

$$I - I^n \approx \frac{I^h - I^{2h}}{2^k - 1}. \quad (13.31)$$

Для формул прямоугольников и трапеций $k = 2$ (см. (13.25) и (13.26)), а для формулы Симпсона $k = 4$ (см. (13.27)); поэтому для этих квадратурных формул равенство (13.31) принимает следующий вид:

$$I - I_{\text{пр}}^h \approx \frac{1}{3} (I_{\text{пр}}^h - I_{\text{пр}}^{2h}), \quad (13.32)$$

$$I - I_{\text{пп}}^h \approx \frac{1}{3} (I_{\text{пп}}^h - I_{\text{пп}}^{2h}), \quad (13.33)$$

$$I - I_C^h \approx \frac{1}{15} (I_C^h - I_C^{2h}). \quad (13.34)$$

Пример 13.5. Применив правило Рунге, оценим погрешность приближенных значений $I_{\text{пп}}^h = 0.74713088$, $I_{\text{пп}}^{2h} = 0.74621079$, $I_C^h = 0.74628418$, полученных в примере 13.2 при вычислении интеграла

$$I = \int_0^1 e^{-x^2} dx \text{ и использующих формулы прямоугольников, трапеций}$$

и Симпсона с шагом $h = 0.1$.

Вычислим приближенные значения интеграла по указанным квадратурным формулам с удвоенным значением шага. В результате получим $I_{\text{пп}}^{2h} = 0.74805326$, $I_{\text{пп}}^{2h} = 0.74436832$, $I_C^{2h} = 0.74682495$.

Применяя теперь формулы (13.32)–(13.34), находим

$$I - I_{\text{пп}}^h \approx \frac{1}{3} (0.74713088 - 0.74805326) \approx -3 \cdot 10^{-4},$$

$$I - I_{\text{пп}}^h \approx \frac{1}{3} (0.74621079 - 0.74436832) \approx 6 \cdot 10^{-4},$$

$$I - I_C^h \approx \frac{1}{15} (0.74682418 - 0.74682495) \approx -5 \cdot 10^{-8}. \blacksquare$$

Естественно, что кроме правила Рунге существуют и другие способы апостериорной оценки погрешности. Например, можно использовать значения $I_{\text{пп}}^h$ и $I_{\text{пп}}^{2h}$, вычисленные по формулам прямоугольников и трапеций с одним и тем же шагом, для практической оценки погрешности каждого из этих значений. Действительно, в равенствах (13.25), (13.26) $C_{\text{тр}} = -2C_{\text{пп}}$, поэтому

$$I_{\text{пп}}^h - I_{\text{пп}}^{2h} = (C_{\text{пп}} - C_{\text{тр}})h^2 + o(h^2) \approx 3C_{\text{пп}}h^2.$$

Откуда следует, что

$$I - I_{\text{пп}}^h \approx \frac{1}{3} (I_{\text{пп}}^h - I_{\text{пп}}^{2h}), \quad (13.35)$$

$$I - I_{\text{пп}}^h \approx -\frac{2}{3} (I_{\text{пп}}^h - I_{\text{пп}}^{2h}). \quad (13.36)$$

Пример 13.6. Применив формулы (13.35), (13.36), оценим погрешности значений $I_{\text{пр}}^h = 0.74713088$, $I_{\text{тр}}^h = 0.74621079$, являющихся приближениями к значению интеграла $I = \int_0^1 e^{-x^2} dx$ (см. пример 13.2).

Имеем $I - I_{\text{пр}}^h \approx \frac{1}{3} (0.74621079 - 0.74713088) \approx -3 \cdot 10^{-4}$, $I - I_{\text{тр}}^h \approx -2(-3 \cdot 10^{-4}) = 6 \cdot 10^{-4}$. Отметим, что полученные оценки совпадают с соответствующими оценками из примера 13.5. ▲

Наличие некоторого правила получения апостериорной оценки погрешности позволяет строить процедуры вычисления интеграла I с заданной точностью ε , достижаемой последовательным дроблением шага интегрирования. Простейшая процедура такого типа состоит в последовательном вычислении значений I^{h_i} и соответствующих апостериорных оценок погрешности ε_i (например, по правилу Рунге) для $h_i = h_0 / 2^i$, где h_0 — начальное значение шага, $i = 1, 2, \dots$. Вычисления прекращаются тогда, когда при некотором i оказывается $|\varepsilon_i| < \varepsilon$ (требуемая точность достигнута) либо тогда, когда величина $|\varepsilon_i|$ начинает возрастать (точность не может быть достигнута из-за влияния вычислительной погрешности).

Пример 13.7. Найдем значение интеграла $I = \int_0^1 e^{-x^2} dx$ с точностью $\varepsilon = 10^{-4}$, используя формулу трапеций и применив процедуру последовательного дробления шага интегрирования, описанную выше.

Возьмем $h_0 = 0.2$. Значения $I^{h_0} = 0.74436832$, $I^{h_1} = 0.74621079$ (где $h_1 = h_0/2 = 0.1$) и $\varepsilon_1 = 6 \cdot 10^{-4}$ были уже получены (см. пример 13.5). Так как $|\varepsilon_1| > \varepsilon$, то уменьшаем шаг вдвое: $h_2 = h_1/2 = 0.05$ и вычисляем $I^{h_2} = 0.74667084$, $\varepsilon_2 = \frac{1}{3} (I^{h_2} - I^{h_1}) = \frac{1}{3} (0.74667084 - 0.74621079) \approx 1.5 \cdot 10^{-4}$.

Так как $|\varepsilon_2| > \varepsilon$, то снова дробим шаг: $h_3 = h_2/2 = 0.025$, вычисляем $I^{h_3} = 0.74678581$, $\varepsilon_3 = \frac{1}{3} (I^{h_3} - I^{h_2}) = \frac{1}{3} (0.74678581 - 0.74667084) \approx 4 \cdot 10^{-5}$. Поскольку $|\varepsilon_3| < \varepsilon$, требуемая точность достигнута и с учетом округления получаем $I = 0.7468 \pm 0.0001$. ▲

3. Экстраполяция Ричардсона. Приближенное равенство (13.30) позволяет получить уточненное значение интеграла

$$I \approx I^{h/2} + \frac{1}{2^k - 1} (I^{h/2} - I^h). \quad (13.37)$$

Таким образом, квадратурная формула I^h порождает новую квадратурную формулу (13.37), имеющую более высокий порядок точности. Если этот порядок известен, то процесс уточнения можно продолжить.

Предположим, например, что для погрешности квадратурной формулы справедливо представление

$$I - I^h = C_1 h^{k_1} + C_2 h^{k_2} + \dots + C_N h^{k_N} + o(h^{k_N}) \quad (13.38)$$

при всех $N = 1, 2, \dots$, причем $0 < k_1 < k_2 < \dots < k_N < \dots$. В этом случае формула (13.37) приводит к следующему методу уточнения, который называют также *методом экстраполяции Ричардсона*¹. Сначала полагают $I_0^h = I^h$. Для вычисления всех последующих приближений используют рекуррентное соотношение

$$I_N^h = I_{N-1}^{h/2} + \frac{1}{2^{k_N} - 1} (I_{N-1}^{h/2} - I_{N-1}^h), \quad N = 1, 2, \dots$$

Для погрешности составной формулы трапеций с постоянным шагом h представление (13.38) действительно имеет место, причем $k_1 = 2$, $k_2 = 4$, ..., $k_N = 2N$. Приведем соответствующий результат.

Теорема 13.5. Пусть функция f является $2N + 1$ раз непрерывно дифференцируемой на отрезке $[a, b]$. Тогда справедлива формула Эйлера—Маклорена

$$\begin{aligned} I - I_{\text{тр}}^h &= -\frac{B_2}{2!} [f'(b) - f'(a)]h^2 - \frac{B_4}{4!} [f^{(3)}(b) - f^{(3)}(a)]h^4 - \dots \\ &\dots - \frac{B_{2N}}{(2N)!} [f^{(2N-1)}(b) - f^{(2N-1)}(a)]h^{2N} + O(h^{2N+1}). \end{aligned}$$

Здесь $B_{2j} = (-1)^{j-1}(2j)! \sum_{k=1}^{\infty} \frac{2}{(2\pi k)^{2j}}$ ($j=1, 2, \dots, N$) — числа Бернуlli. ■

Таким образом, к формуле трапеции можно применить экстраполяцию Ричардсона. В результате получается так называемый *метод Ромберга*. Существуют стандартные программы вычисления интегралов методом Ромберга. Правда, следует отметить, что эффективность этого метода не слишком велика. Как правило, лучший результат дает приме-

¹ Арчибалд Рид Ричардсон (1881—1954) — английский математик.

нение квадратурных формул Гаусса или рассматриваемых ниже адаптивных процедур численного интегрирования.

Замечание. Первый же шаг метода Ромберга приводит к уточнению квадратурной формулы трапеций, совпадающему с формулой Симпсона. Действительно,

$$\begin{aligned} I_{\text{тр}}^{h/2} + \frac{1}{3} (I_{\text{тр}}^{h/2} - I_{\text{тр}}^h) &= \frac{4}{3} I_{\text{тр}}^{h/2} - \frac{1}{3} I_{\text{тр}}^h = \frac{2h}{3} \left[\frac{f_0 + f_n}{2} + \sum_{i=1}^{2n-1} f_{i/2} \right] - \\ &- \frac{h}{3} \left[\frac{f_0 + f_n}{2} + \sum_{i=1}^{n-1} f_i \right] = \frac{h}{6} \left(f_0 + f_n + 4 \sum_{i=1}^n f_{i-1/2} + 2 \sum_{i=1}^{n-1} f_i \right) = I_C^h. \end{aligned}$$

4. Адаптивные процедуры численного интегрирования. До сих пор нам было удобнее рассматривать методы численного интегрирования с постоянным шагом. Однако они обладают значительно меньшей эффективностью и используются в вычислительной практике существенно реже, чем методы с переменным шагом. Объясняется это тем, что, равномерно распределяя узлы по отрезку интегрирования, мы полностью игнорируем особенности поведения подынтегральной функции. В то же время интуитивно ясно, что на участках плавного изменения функции достаточно поместить сравнительно небольшое число узлов, разместив значительно большее их число на участках резкого изменения функции. Распределение узлов интегрирования в соответствии с характером поведения подынтегральной функции часто позволяет при том же общем числе узлов получить значительно более высокую точность. Тем не менее выбор соответствующего неравномерного распределения узлов интегрирования является очень сложной задачей и вряд ли мог быть широко использован на практике, если бы для решения этой задачи не удалось привлечь компьютер.

Современные процедуры численного интегрирования (*адаптивные квадратурные программы*) используют некоторый алгоритм автоматического распределения узлов интегрирования. Пользователь такой программы задает отрезок $[a, b]$, правило вычисления функции f и требуемую точность $\varepsilon > 0$. Программа стремится, используя по возможности минимальное число узлов интегрирования, распределять их так, чтобы

найденное значение $I^h = \sum_{i=1}^n I_i^{h_i}$ удовлетворяло неравенству

$$\left| I^h - \int_a^b f(x) dx \right| < \varepsilon. \quad (13.39)$$

Здесь $I_i^{h_i}$ — приближенное значение интеграла $I_i = \int_{x_{i-1}}^{x_i} f(x) dx$, вычисляемое по некоторой формуле.

Типичная программа разбивает исходный отрезок $[a, b]$ на элементарные отрезки $[x_{i-1}, x_i]$ так, чтобы для погрешностей $R_i = I_i - I_i^{h_i}$ выполнялось неравенство

$$\sum_{i=1}^n |R_i| < \varepsilon. \quad (13.40)$$

При этом каждый из элементарных отрезков получается, как правило, делением пополам одного из отрезков, найденных на более раннем шаге алгоритма. Заметим, что неравенство (13.40) выполняется, если каждая из погрешностей R_i удовлетворяет условию $|R_i| < h_i \varepsilon / (b - a)$. Действительно, тогда

$$|I - I^h| = \left| \sum_{i=1}^n (I_i - I_i^{h_i}) \right| \leq \sum_{i=1}^n |R_i| < \frac{\varepsilon}{b-a} \sum_{i=1}^n h_i = \varepsilon.$$

Рассмотрим, например, одну из простейших адаптивных процедур, основанную на формуле трапеций

$$I_i \approx I_{\text{tp}, i}^{h_i} = \frac{h_i}{2} (f_{i-1} + f_i) \quad (13.41)$$

и использующую для контроля точности составную формулу трапеций с шагом $h_i/2$:

$$I_i \approx I_i^{h_i} = \frac{h_i}{4} (f_{i-1} + 2f_{i-1/2} + f_i). \quad (13.42)$$

Заметим, что если значение $I_{\text{tp}, i}^{h_i}$, уже найдено, то для нахождения значения $I_i^{h_i}$ требуется лишь одно дополнительное вычисление функции f в точке $x_{i-1/2}$.

Можно показать, что для оценки погрешностей квадратурных формул интерполяционного типа на элементарных отрезках правило Рунге сохраняет силу, поэтому погрешность приближенного значения $I_i^{h_i}$ можно оценить по формуле

$$R_i \approx \varepsilon_i = \frac{1}{3} (I_i^{h_i} - I_{\text{tp}, i}^{h_i}). \quad (13.43)$$

В рассматриваемой аддитивной процедуре последовательно выбирают точки x_i и вычисляют значения

$$S_i = \sum_{j=1}^i I_j^{h_j} \approx \int_a^{x_i} f(x) dx \quad (i=1, 2, \dots, n),$$

последнее из которых S_n совпадает с I^h и принимается за приближенное значение интеграла I . Перед началом работы полагают $S_0 = 0$, $x_0 = a$ и задают некоторое начальное значение шага h_1 .

Опишем i -й шаг процедуры в предположении, что значение S_{i-1} уже найдено и очередное значение шага h_i определено.

1°. По формулам (13.41)–(13.43) вычисляют значения $I_{\text{тр}, i}^{h_i}$, $I_i^{h_i}$, ε_i .

2°. Если $|\varepsilon_i| > h_i \varepsilon / (b - a)$, то шаг h_i уменьшают в 2 раза и повторяются вычисления п. 1°.

3°. После того как очередное дробление шага приводит к выполнению условия $|\varepsilon_i| \leq h_i \varepsilon / (b - a)$, вычисляют значения $x_i = x_{i-1} + h_i$, $S_i = S_{i-1} + I_i^{h_i}$.

4°. Если $x_i + h_i > b$, то полагают $h_{i+1} = b - x_i$. В противном случае полагают $h_{i+1} = h_i$. На этом i -й шаг завершается.

Замечание 1. В некоторых аддитивных процедурах при выполнении условия типа $|\varepsilon_i| \ll h_i \varepsilon / (b - a)$ очередное значение шага удваивается: $h_{i+1} = 2h_i$. Таким образом, шаг интегрирования в зависимости от характера поведения подынтегральной функции может не только измельчаться, но и укрупняться.

Замечание 2. Известно [9], что при оптимальном распределении узлов интегрирования модули погрешностей, приходящихся на элементарные отрезки интегрирования, должны быть примерно одинаковыми. Распределение узлов, полученное с помощью аддитивных процедур, как правило, не является таковым. Однако для большинства функций оно является вполне удовлетворительным.

5. Квадратурные формулы Гаусса—Кронрода. Известно, что квадратурные формулы Гаусса с разным числом узлов не имеют общих узлов, поэтому попытка использовать для апостериорной оценки погрешности формулы Гаусса G_{N+1} с $N+1$ узлами формулу G_N с N узлами приводит к необходимости вычисления $2N+1$ значений подынтегральной функции.

В 1965 г. советский математик А.С. Кронрод предложил в дополнение формулы Гаусса G_N использовать специальную квадратурную формулу K_{2N+1} с $2N+1$ узлами, которая использует только $N+1$ новый узел. Таким образом, вычисление пары Гаусса—Кронрода (G_N, K_{2N+1}) требует такой же вычислительной работы, как и вычисление пары (G_N, G_{N+1}). Однако формула K_{2N+1} точна для многочленов степени $3N+1$, а формула G_{N+1} — лишь для многочленов степени $2N+1$.

Использование пары Гаусса—Кронрода (G_N, K_{2N+1}) в сочетании с оценкой погрешности вида

$$|I - K_{2N+1}| \approx (200 |G_N - K_{2N+1}|)^{3/2}$$

приводит к одной из самых эффективных в настоящее время адаптивных процедур численного интегрирования.

Результаты А.С. Кронрода опубликованы в книге [59]. Подробнее о квадратурных формулах Гаусса—Кронрода см. в [54].

§ 13.5. Вычисление интегралов в нерегулярных случаях

Нередко приходится вычислять интегралы

$$\int_a^b F(x) dx \quad (13.44)$$

от функций, имеющих те или иные особенности. Например, сама функция F (или ее производная некоторого порядка) имеет участки резкого изменения, точки разрыва или является неограниченной. Такие функции плохо аппроксимируются многочленами и поэтому для вычисления соответствующих интегралов может оказаться неэффективным непосредственное применение стандартных квадратурных формул, рассчитанных на возможность кусочно-многочленной аппроксимации функции F . Случай, когда в интеграле (13.44) промежуток интегрирования бесконечен, также требует специального рассмотрения.

Пример 13.8. Функция $\frac{1}{\sqrt{x}} e^{-x^2}$ имеет особенность в точке $x = 0$.

Попробуем применить для вычисления интеграла

$$I = \int_0^1 \frac{1}{\sqrt{x}} e^{-x^2} dx \approx 1.689677 \quad (13.45)$$

формулу прямоугольников с постоянным шагом

$$I \approx I_{\text{пп}}^h = \sum_{i=1}^n \frac{1}{\sqrt{x_i - 1/2}} e^{-x_i^2/2} h. \quad (13.46)$$

Результаты вычислений для нескольких значений шага h приведены в табл. 13.3. Заметим, что значения $I_{\text{пп}}^h$ сходятся к I очень медленно.

Таблица 13.3

h	$I_{\text{пп}}^h$	$I - I_{\text{пп}}^h$
0.200	1.420	$2.7 \cdot 10^{-1}$
0.100	1.499	$1.9 \cdot 10^{-1}$
0.050	1.555	$1.4 \cdot 10^{-1}$
0.025	1.594	$9.6 \cdot 10^{-2}$

Теоретический анализ погрешности показывает, что она убывает пропорционально лишь $h^{1/2}$, а не h^2 как в регулярном случае. ▲

Укажем на некоторые подходы к вычислению интегралов в нерегулярных случаях, позволяющие учесть особенности поведения функции F и благодаря этому значительно сократить затраты машинного времени или достигнуть большей точности.

1. Разбиение промежутка интегрирования на части. Пусть подынтегральная функция F является кусочно-гладкой и $c_1 < c_2 < \dots < c_p$ — известные точки разрыва функции F либо ее производных. В этом случае имеет смысл представить интеграл (13.44) в виде суммы:

$$\int_a^b F(x) dx = \int_a^{c_1} F(x) dx + \int_{c_1}^{c_2} F(x) dx + \dots + \int_{c_{p-1}}^b F(x) dx. \quad (13.47)$$

Вычисление каждого из входящих в сумму (13.47) интегралов представляет собой стандартную задачу, так как на каждом из частичных отрезков $[a, c_1], [c_1, c_2], \dots, [c_p, b]$ подынтегральная функция является гладкой.

Разбиение промежутка интегрирования на части может оказаться полезным приемом и в других ситуациях, например тогда, когда имеет смысл на разных частях применять различные квадратурные формулы. Другой пример дает стандартный прием вычисления несобственных

интегралов вида $\int_a^{\infty} F(x) dx$. Если требуется вычислить такой интеграл с точностью $\epsilon > 0$, то его представляют в виде суммы:

$$\int_a^{\infty} F(x) dx = \int_a^b F(x) dx + \int_b^{\infty} F(x) dx.$$

Затем благодаря выбору достаточно большого значения b добиваются выполнения неравенства $|\int_b^{\infty} F(x) dx| < \epsilon/2$ и вычисляют интеграл $\int_a^b F(x) dx$ с точностью $\epsilon/2$.

2. Выделение веса. В некоторых случаях подынтегральная функция допускает разложение на два сомножителя $F(x) = \rho(x) \cdot f(x)$, где функция $\rho(x)$ является достаточно простой и имеет те же особенности, что и $F(x)$, а $f(x)$ — гладкая функция. Тогда имеет смысл рассматривать интеграл (13.44) в виде

$$I = \int_a^b \rho(x) f(x) dx. \quad (13.48)$$

Здесь функция $\rho(x)$ называется *весовой функцией* (или *весом*). При построении численных методов вычисления интеграла (13.48) весовая функция считается фиксированной. В то же время $f(x)$ может быть произвольной достаточно гладкой функцией.

Примерами весовых функций могут служить постоянный вес $\rho(x) \equiv 1$, весовые функции Якоби $\rho(x) = (x - a)^{\alpha} (b - x)^{\beta}$ ($a < x < b$, $-1 < \alpha, -1 < \beta$), Лагерра¹ $\rho(x) = x^{\alpha} e^{-x}$ ($0 < x < \infty$) и Эрмита $\rho(x) = e^{-x^2}$ ($-\infty < x < \infty$),

соответствующие интегралам вида $\int_a^b f(x) dx$, $\int_a^b (x - a)^{\alpha} (b - x)^{\beta} f(x) dx$,

$$\int_0^{\infty} x^{\alpha} e^{-x} f(x) dx, \quad \int_{-\infty}^{\infty} e^{-x^2} f(x) dx.$$

Методы приближенного вычисления интегралов, рассмотренные в предыдущих параграфах, применимы и к задаче вычисления интегралов с весом.

¹ Эдмунд Никола Лагерр (1834—1886) — французский математик.

Пусть $P_{m,i}$ — интерполяционные многочлены (13.17). Приближенная замена интеграла (13.48) суммой:

$$I^h = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} \rho(x) P_{m-i}(x) dx$$

приводит к следующей квадратурной формуле интерполяционного типа

$$I \approx I^h = \sum_{i=1}^n h_i \sum_{j=0}^m a_{ij} f(x_{i-1/2} + t_j h_i / 2). \quad (13.49)$$

Здесь коэффициенты a_{ij} вычисляются по формуле

$$a_{ij} = \frac{1}{2} \int_{-1}^1 \rho(x_{i-1/2} + th_i / 2) \prod_{\substack{k=0 \\ k \neq j}}^m \frac{t - t_k}{t_j - t_k} dt.$$

Если функция f имеет на отрезке $[a, b]$ непрерывную производную порядка $m+1$, то для погрешности формулы (13.49) верна оценка

$$|I - I^h| \leq C_m M_{m+1} \int_a^b |\rho(x)| dx \cdot h_{\max}^{m+1}.$$

Пример 13.9. Выведем аналог формулы прямоугольников с постоянным шагом для вычисления интеграла

$$I = \int_0^1 \frac{1}{\sqrt{x}} f(x) dx \quad (13.50)$$

Заменяя функцию f на элементарном отрезке $[x_{i-1}, x_i]$ постоянной $f_{i-1/2}$ и учитывая, что

$$\int_{x_{i-1}}^{x_i} \frac{1}{\sqrt{x}} f(x) dx \approx \int_{x_{i-1}}^{x_i} \frac{1}{\sqrt{x}} f_{i-1/2} dx = \frac{2h}{\sqrt{x_{i-1}} + \sqrt{x_i}} f_{i-1/2},$$

получаем следующую квадратурную формулу:

$$\int_0^1 \frac{1}{\sqrt{x}} f(x) dx \approx I^h = h \sum_{i=1}^n \frac{2}{\sqrt{x_{i-1}} + \sqrt{x_i}} f_{i-1/2}. \quad \blacktriangle \quad (13.51)$$

Пример 13.10. Применим квадратурную формулу (13.51) для вычисления интеграла (13.45) при тех же значениях шага, что и в примере 13.9.

В рассматриваемом случае формула (13.51) принимает вид

$$\int_0^1 \frac{1}{\sqrt{x}} e^{-x^2} dx \approx h \sum_{i=1}^n \frac{2}{\sqrt{x_{i-1}} + \sqrt{x_i}} e^{-x_i^2 - 1/2}. \quad (13.52)$$

Полученные с ее помощью результаты приведены в табл. 13.4.

Таблица 13.4

h	I_h	$I - I^h$
0.200	1.686276	$3.4 \cdot 10^{-3}$
0.100	1.688958	$7.2 \cdot 10^{-4}$
0.050	1.689521	$1.6 \cdot 10^{-4}$
0.025	1.689642	$3.5 \cdot 10^{-5}$

Сравнение с результатами примера 13.9 показывает, что для вычисления интеграла (13.45) формула (13.52) имеет безусловное преимущество перед формулой прямоугольников (13.46). ▲

Для вычисления интегралов (13.48) применяют и квадратурные формулы Гаусса $\int_a^b p(x)f(x) dx \approx \sum_{i=1}^n A_i f(x_i)$, точные для многочленов наиболее высокой степени. Они строятся аналогично тому, как это было сделано для постоянного веса $p(x) \equiv 1$ (см. § 13.4).

Пример 13.11. Для вычисления интеграла (13.50) построим квадратурную формулу Гаусса с одним узлом

$$\int_0^1 \frac{1}{\sqrt{x}} f(x) dx \approx c_1 f(x_1). \quad (13.53)$$

Потребуем, чтобы формула (13.53) была точна для многочленов первой степени. Это эквивалентно выполнению равенств $2 = \int_0^1 \frac{1}{\sqrt{x}} dx = c_1$,

$\frac{2}{3} = \int_0^1 \frac{x}{\sqrt{x}} dx = c_1 x_1$. Таким образом $c_1 = 2$, $x_1 = 1/3$ и формула (13.53)

принимает вид

$$\int_0^1 \frac{1}{\sqrt{x}} f(x) dx \approx 2f\left(\frac{1}{3}\right). \quad (13.54)$$

Хотя формула (13.54) и кажется примитивной, применяя ее для вычисления интеграла (13.45), получаем значение $I \approx 1.789679$, абсолютная погрешность которого равна 0.1 и практически совпадает с погрешностью значения, найденного в примере 13.9 по формуле прямоугольников с шагом $h = 0.025$. Замечательно то, что формула Гаусса (13.54) достигает точности $\epsilon = 0.1$ при использовании только одного вычисления значения функции $f(x) = e^{-x^2}$. В то же время формула (13.46) для достижения той же точности требует вычисления 40 значений функции. ▲

3. Формула Эрмита. Для вычисления интегралов вида $\int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}} dx$,

т.е. при $\rho(x) = 1/\sqrt{1-x^2}$, $a = -1$, $b = 1$, используют квадратурную формулу Гаусса

$$\int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}} dx \approx \frac{\pi}{n} \sum_{i=1}^n f(x_i),$$

называемую *формулой Эрмита*. Узлами этой формулы являются нули многочлена Чебышева $T_n(x)$, т.е. числа $x_i = \cos \frac{(2i-1)\pi}{2n}$, $1 \leq i \leq n$.

4. Интегрирование быстро осциллирующих функций. В задачах радиотехники часто встречается проблема вычисления интегралов вида¹

$$\int_a^b f(x) e^{i\omega x} dx = \int_a^b f(x) \cos \omega x dx + i \int_a^b f(x) \sin \omega x dx. \quad (13.55)$$

Здесь $\omega(b-a) \gg 1$, $f(x)$ — некоторая достаточно гладкая функция, а i — мнимая единица. Функции $f(x) \cos \omega x$ и $f(x) \sin \omega x$ являются быстро меняющимися и имеют на отрезке $[a, b]$ порядка $(b-a)\omega/\pi$ нулей. Если попытаться вычислить интеграл (13.55) с помощью стандартных квадратурных формул, то для обеспечения приемлемой точности на каждый «полупериод» колебаний подынтегральной функции потребуется поместить хотя бы несколько (например, порядка десяти) точек. Так как на отрезок $[a, b]$ приходится примерно $(b-a)\omega/\pi$ таких «полупериодов», то необходимо по меньшей мере порядка $\omega(b-a) \gg 1$ узлов интегрирова-

¹ Например, $F(x) = f(x) e^{i\omega x}$ отвечает несущему высокочастотному колебанию $e^{i\omega x}$ с модулированной амплитудой $f(x)$.

ния. Следовательно, стандартный подход к вычислению интегралов вида (13.55) потребует слишком больших затрат машинного времени.

Для существенного уменьшения объема вычислений в равенстве (13.55) полезно рассматривать функцию $\rho(x) = e^{i\omega x}$ как весовую. Тогда кусочно-полиномиальная интерполяция функции $f(x)$ приводит к квадратурным формулам интерполяционного типа, которые принято называть *формулами Филона*.

Выведем одну из таких формул, основанную на интерполяции функции $f(x)$ на каждом из элементарных отрезков $[x_{k-1}, x_k]$ линейной функцией $P_{1,k}(x) = f_{k-1} + (f_k - f_{k-1})(x - x_{k-1})/h_k$ и являющуюся аналогом составной квадратурной формулы трапеций. Положим

$$I_k = \int_{x_{k-1}}^{x_k} f(x) e^{i\omega x} dx \approx I_k^{h_k} = \int_{x_{k-1}}^{x_k} P_{1,k}(x) e^{i\omega x} dx.$$

Вычислив интеграл $I_k^{h_k}$, получим формулу

$$I_k^{h_k} = \frac{h_k}{2} (A_k f_{k-1} + B_k f_k) e^{i\omega x_{k-1}/2}.$$

Здесь

$$A_k = \frac{\sin p_k}{p_k} + i \frac{p_k \cos p_k - \sin p_k}{p_k^2}, \quad B_k = \frac{\sin p_k}{p_k} - i \frac{p_k \cos p_k - \sin p_k}{p_k^2}, \quad p_k = \frac{\omega h_k}{2}.$$

В результате приходим к составной формуле вида

$$\int_a^b f(x) e^{i\omega x} dx \approx \sum_{k=1}^n \frac{h_k}{2} (A_k f_{k-1} + B_k f_k) e^{i\omega x_{k-1}/2}.$$

5. Аддитивное выделение особенности. Иногда подынтегральную функцию удается представить в виде суммы $F(x) = \phi(x) + \psi(x)$, где функция $\phi(x)$ содержит особенность, но интегрируется аналитически, а функция $\psi(x)$ является достаточно гладкой. Тогда интеграл от функции F представляют в виде суммы двух интегралов:

$$\int_a^b F(x) dx = \int_a^b \phi(x) dx + \int_a^b \psi(x) dx = I^{(1)} + I^{(2)}.$$

Первый из них вычисляется аналитически, а значение второго можно найти с помощью той или иной квадратурной формулы.

Пример 13.12. Указанный прием можно использовать для вычисления интеграла (13.45). Представим интеграл в виде

$$\int_0^1 \frac{1}{\sqrt{x}} e^{-x^2} dx = \int_0^1 \frac{1-x^2}{\sqrt{x}} dx + \int_0^1 \frac{e^{-x^2}-1+x^2}{\sqrt{x}} dx = I^{(1)} + I^{(2)}.$$

Интеграл $I^{(1)}$ вычисляется аналитически: $I^{(1)} = 1.6$. В то же время функция $(e^{-x^2} - 1 + x^2)/\sqrt{x}$ трижды непрерывно дифференцируема на отрезке $[0, 1]$. Поэтому интеграл $I^{(2)}$ можно вычислить по формуле прямоугольников. В результате приходим к формуле

$$I \approx I^h = 1.6 + h \sum_{i=1}^n \frac{1}{\sqrt{x_{i-1/2}}} \left(e^{-x_{i-1/2}^2} - 1 + x_{i-1/2}^2 \right).$$

Найденные по ней приближенные значения интеграла приведены в табл. 13.5. ▲

Таблица 13.5

h	I^h	$I - I^h$
0.200	1.687874	$1.8 \cdot 10^{-3}$
0.100	1.689227	$4.5 \cdot 10^{-4}$
0.050	1.689565	$1.1 \cdot 10^{-4}$
0.025	1.689649	$2.8 \cdot 10^{-5}$

6. Замена переменой и сгущение сетки. Конечно, отмеченные приемы представляют лишь небольшую часть тех средств, которые применяются при вычислении интегралов в нерегулярных случаях. Иногда, например, оказывается полезной замена переменой $x = \phi(t)$, приводящая интеграл к виду

$$\int_a^b F(x) dx = \int_a^\beta F(\phi(t))\phi'(t) dt.$$

Если функция $g(t) = F(\phi(t))\phi'(t)$ достаточно гладкая, то интеграл может уже быть вычислен стандартными методами

Пример 13.13. Сделаем в интеграле (13.45) замену $x = t^2$. Тогда

$$I = \int_0^1 \frac{1}{\sqrt{x}} e^{-x^2} dx = 2 \int_0^1 e^{-t^4} dt.$$

Для вычисления полученного интеграла можно уже использовать одну из стандартных квадратурных формул, например, — составную формулу прямоугольников с узлами $t_i = i/n$, $0 \leq i \leq n$ и шагом $h = 1/n$

$$I = 2 \int_0^1 e^{-t^4} dt \approx 2h \sum_{i=1}^n e^{-t_{i-1/2}^4}.$$

Здесь $t_{i-1/2} = (i - 1/2)/n$, $1 \leq i \leq n$. ▲

Близкий прием состоит в сгущении сетки узлов x_i вблизи особенности.

Пример 13.14. Например, для вычисления интеграла (13.45) можно использовать составную квадратурную формулу прямоугольников

$$I \approx \sum_{i=1}^n \frac{1}{\sqrt{x_{i-1/2}}} e^{-x_{i-1/2}^2} h_i$$

с узлами $x_i = (i/n)^2$, $0 \leq i \leq n$, сгущающимися вблизи точки $x = 0$. Здесь $x_{i-1/2} = (x_{i-1} + x_i)/2$, $h_i = x_i - x_{i-1}$, $1 \leq i \leq n$. ▲

Иногда заменой переменной $x = \phi(t)$ удается свести проблему вычисления несобственного интеграла по бесконечном интервалу к интегралу по конечному отрезку с подынтегральной функцией, не имеющей особенностей.

Пример 13.5. Заменой переменной $x = 1/t$ интеграл $\int_1^\infty \frac{\sqrt{1+x^2}}{\sqrt{1+x^6}} dx$

преобразуется в равный ему интеграл $\int_0^1 \frac{\sqrt{t^2+1}}{\sqrt{t^6+1}} dt$, для вычисления

которого можно использовать стандартные квадратурные формулы. ▲

§ 13.6. Дополнительные замечания

1. Мы не рассматриваем проблему вычисления кратного интеграла

$$I = \iint_G \dots \int f(x_1, x_2, \dots, x_m) dx_1 dx_2 \dots dx_m, \quad (13.56)$$

поскольку это потребовало бы привлечения достаточно сложного математического аппарата. Ограничимся указанием на то, что в принципе вычисление интегралов (13.56) можно проводить методами, аналогич-

ными рассмотренным в этой главе. Соответствующие *кубатурные формулы*¹ для вычисления кратных интегралов имеют вид

$$I \approx \sum_{j=1}^N A_j f(x_1^{(j)}, x_2^{(j)}, \dots, x_m^{(j)}) \quad (13.57)$$

Среди формул (13.57) есть кубатурные формулы интерполяционного типа и кубатурные формулы Гаусса. Иногда для вычисления кратного интеграла оказывается целесообразным сведение его к повторному вычислению однократных интегралов. Для первоначального знакомства с методами вычисления кратных интегралов можно рекомендовать книги [9, 51].

2. Вычисление кратных интегралов уже при не очень больших значениях $m \geq 6$ является очень сложной задачей. Применение для вычисления таких интегралов кубатурных формул типа (13.57) требует (даже при очень скромных запросах к точности) такого большого числа N вычислений значений функции f , что решение задачи даже при использовании самых современных компьютеров становится нереальным. Привлекательной альтернативой в такой ситуации становится использование *метода Монте-Карло*. Простейшее представление об этом методе (на примере вычисления однократного интеграла) можно получить из учебника [25]. Мы все же рекомендуем обратиться и к весьма содержательному обсуждению метода Монте-Карло, проведенному в книге [9].

3. Иногда возникает необходимость по известной функции $f(x)$, заданной на отрезке $[a, b]$, восстановить ее первообразную

$$y(x) = \int_a^x f(\xi) d\xi, \quad a \leq x \leq b. \quad (13.58)$$

При каждом фиксированном x функцию (13.58) можно рассматривать как определенный интеграл вида (13.1) и вычислять с помощью одного из известных методов. Однако если требуется находить значения $y(x)$ в большом числе различных точек, то такой подход становится нецелесообразным. Оказывается более выгодным разбить отрезок $[a, b]$ на элементарные отрезки точками $a = x_0 < x_1 < x_2 < \dots < x_n = b$, а затем составить таблицу значений $y_i \approx y(x_i)$, $0 \leq i \leq n$. Значения y_i можно найти, например по формуле $y_i = y_{i-1} + I_i^{h_i}$, $1 \leq i \leq n$. Здесь $y_0 = 0$, а $I_i^{h_i}$ — при-

¹ Впрочем формулы (13.57) называют также и квадратурными.

ближение к интегралу $\int_{x_{i-1}}^{x_i} f(x) dx$, полученное с помощью одной из квадратурных формул. Значение $y(x)$ в любой из промежуточных точек можно затем приближенно восстановить, использовав интерполяцию. Так как значения $y'(x_i) = f(x_i)$ фактически также известны, то весьма подходящим для интерполяции на каждом элементарном отрезке $[x_{i-1}, x_i]$ является кубический многочлен Эрмита (см. § 11.5). Использование этого способа интерполяции позволяет находить значения $y(x)$ с довольно высокой точностью по сравнительно редкой таблице значений.

4. В данной главе в основном обсуждались не вычислительные алгоритмы, а методы дискретизации, т.е. методы замены определенных интегралов соответствующими квадратурными суммами. Как бы ни был организован алгоритм, он все же предполагает вычисление квадратурной суммы. С увеличением числа слагаемых возрастает влияние вычислительной погрешности на результат суммирования. При очень больших значениях N даже для хорошо обусловленных квадратурных формул соответствующий вычислительный алгоритм может стать плохо обусловленным. Тем не менее при умеренном значении числа узлов влияние погрешностей округления невелико и им часто можно пренебречь.

Глава 14

ЧИСЛЕННЫЕ МЕТОДЫ РЕШЕНИЯ ЗАДАЧИ КОШИ ДЛЯ ОБЫКНОВЕННЫХ ДИФФЕРЕНЦИАЛЬНЫХ УРАВНЕНИЙ

Часто приходится иметь дело с процессами, характеристики которых непрерывным образом меняются со временем¹ t . Соответствующие явления, как правило, подчиняются физическим законам, которые формулируются в виде дифференциальных уравнений. Одной из основных математических задач, которые приходится решать для таких уравнений, является *задача Коши* (или *начальная задача*). Чаще всего к ней приходят тогда, когда начальное состояние некоторой физической системы в момент времени t_0 считается известным, и требуется предсказать ее поведение при $t \geq t_0$. Понимание того, что задача Коши описывает развитие тех или иных процессов во времени, значительно упрощает восприятие как подходов к ее решению, так и критериев оценки качества получаемых приближений.

Подавляющее большинство возникающих на практике начальных задач невозможно решить без использования вычислительной техники. Поэтому в прикладных расчетах численные методы решения задачи Коши играют особую роль.

Моделирование самых разнообразных процессов приводит к необходимости решать системы дифференциальных уравнений (иногда довольно высокого порядка). Тем не менее большая часть этой главы (§ 14.1—14.9) посвящена рассмотрению методов решения задачи Коши для одного дифференциального уравнения первого порядка. Это традиционный подход, упрощающий как изложение методов, так и понимание их существа. Переход от решения одного уравнения к решению систем дифференциальных уравнений не вызывает затем серьезных затруднений (по крайней мере формального характера). Некоторые особенности решения задачи Коши для систем уравнений изложены в § 14.10 и 14.11. При этом значительное внимание уделяется проблеме устойчивости численных методов и так называемым жестким задачам.

¹ В роли t может выступать и другая (например, пространственная) переменная. Тем не менее в этой главе переменную t будем называть временем.

§ 14.1. Задача Коши для дифференциального уравнения первого порядка

1. Постановка задачи. Напомним, что *решением обыкновенного дифференциального уравнения первого порядка*

$$y'(t) = f(t, y(t)) \quad (14.1)$$

называется дифференцируемая функция $y(t)$, которая при подстановке в уравнение (14.1) обращает его в тождество. График решения дифференциального уравнения называют *интегральной кривой*. Процесс нахождения решений дифференциального уравнения принято называть *интегрированием* этого уравнения.

Исходя из геометрического смысла производной y' заметим, что уравнение (14.1) задает в каждой точке (t, y) плоскости переменных t, y значение $f(t, y)$ тангенса угла α наклона (к оси Ot) касательной к графику решения, проходящего через эту точку. Величину $k = \operatorname{tg} \alpha = f(t, y)$ далее будем называть *угловым коэффициентом* (рис. 14.1). Если теперь в каждой точке (t, y) задать с помощью некоторого вектора направление касательной, определяемое значением $f(t, y)$, то получится так называемое *поле направлений* (рис. 14.2, а). Таким образом, геометрически задача интегрирования дифференциальных уравнений состоит в нахождении в интегральных кривых, которые в каждой своей точке имеют заданное направление касательной (рис. 14.2, б). Для того, чтобы выделить из семейства решений дифференциального уравнения (14.1) одно конкретное решение, задают *начальное условие*

$$y(t_0) = y_0, \quad (14.2)$$

где t_0 — некоторое фиксированное значение аргумента t , y_0 — величина, называемая *начальным значением*.

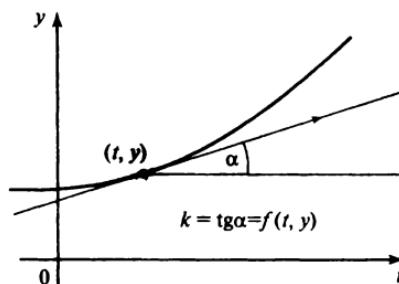
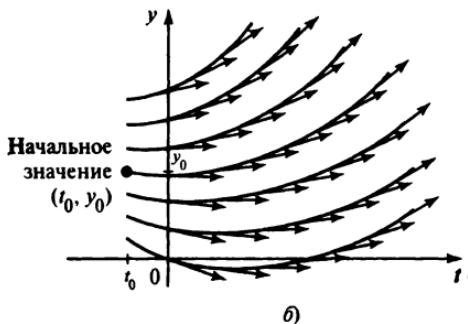
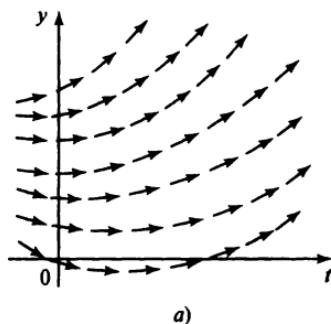


Рис. 14.1



Геометрическая интерпретация использования начального условия состоит в выборе из семейства интегральных кривых той кривой, которая проходит через фиксированную точку (t_0, y_0) .

Задачу нахождения при $t > t_0$ решения $y(t)$ дифференциального уравнения (14.1), удовлетворяющего начальному условию (14.2), будем называть *задачей Коши*. В некоторых случаях представляет интерес поведение решения при всех $t > t_0$. Однако чаще ограничиваются определением решения на конечном отрезке $[t_0, T]$.

2. Разрешимость задачи Коши. Пусть Π_T — множество точек (t, y) , удовлетворяющих условию $t_0 \leq t \leq T, -\infty < y < \infty$; это множество будем называть *полосой*.

Приведем одну из теорем о разрешимости задачи Коши.

Теорема 14.1. Пусть функция $f(t, y)$ определена и непрерывна в полосе Π_T . Предположим также, что она удовлетворяет условию Липшица¹

$$|f(t, y_1) - f(t, y_2)| \leq L |y_1 - y_2| \quad (14.3)$$

для всех $t_0 \leq t \leq T$ и произвольных y_1, y_2 , где L — некоторая постоянная (постоянная Липшица).

Тогда для каждого начального значения y_0 существует единственное решение $y(t)$ задачи Коши (14.1), (14.2), определенное на отрезке $[t_0, T]$. ■

¹ Рудольф Липшиц (1832—1903) — немецкий математик.

Замечание 1. Для дифференцируемых по y функций f условие (14.3) выполняется тогда и только тогда, когда для всех $(t, y) \in \Pi_T$, справедливо неравенство

$$|f'_y(t, y)| \leq L. \quad (14.4)$$

Поэтому условие (14.4) можно также называть *условием Липшица*.

Замечание 2. Теорема 14.1 остается справедливой, если в ее формулировке условие Липшица (14.3) заменить менее ограничительным *односторонним условием Липшица*

$$(f(t, y_1) - f(t, y_2))(y_1 - y_2) \leq \sigma(y_1 - y_2)^2. \quad (14.5)$$

Подчеркнем, что входящая в это условие постоянная σ может иметь произвольный знак.

Для дифференцируемых по y функций f условие (14.5) выполняется тогда и только тогда, когда для всех $(t, y) \in \Pi_T$ справедливо неравенство

$$f'_y(t, y) \leq \sigma. \quad (14.6)$$

Ясно, что для функций, удовлетворяющих условию Липшица с постоянной L , одностороннее условие заведомо выполнено с постоянной $\sigma \leq L$.

Пример 14.1. Функция $f(t) = \cos(t + y)$ удовлетворяет условию Липшица с постоянной $L = 1$, так как $f'_y = -\sin(t + y)$ и $|f'_y| \leq 1$. Отсюда следует, что решение задачи Коши $y' = \cos(t + y)$, $y(t_0) = y_0$ существует и единствено на любом отрезке $[t_0, T]$. ▲

Пример 14.2. Функция $f(t, y) = t - y^3$ не удовлетворяет условию Липшица, поскольку $f'_y = -3y^2$ и модуль этой величины не ограничен. В то же время одностороннее условие (14.6) выполняется с постоянной $\sigma = 0$. Следовательно, можно утверждать, что решение задачи Коши $y' = t - y^3$, $y(t_0) = y_0$ существует и единствено на любом отрезке $[t_0, T]$. ▲

Пример 14.3. Функция $f(t, y) = t + y^3$ не удовлетворяет одностороннему условию (14.6), так как частная производная $f'_y = 3y^2$ не ограничена сверху. Поэтому вопрос о разрешимости задачи Коши $y' = t + y^3$, $y(t_0) = y_0$ требует дополнительного исследования. ▲

Отметим следующий полезный результат, указывающий на зависимость степени гладкости решения задачи Коши от степени гладкости правой части дифференциального уравнения.

Теорема 14.2. Пусть функция f непрерывно дифференцируема m раз в полосе Π_T . Тогда если функция y является на отрезке $[t_0, T]$ решением задачи Коши (14.1), (14.2), то она непрерывно дифференцируема $m + 1$ раз на этом отрезке.

Это утверждение непосредственно вытекает из возможности дифференцирования тождества $y'(t) \equiv f(t, y(t))$ не менее чем m раз. ■

В дальнейшем функции f и y будем предполагать дифференцируемыми столько раз, сколько потребуется при рассмотрении соответствующих численных методов.

3. Устойчивость решения задачи Коши на конечном отрезке. Этот вопрос весьма важен для понимания особенностей методов численного интегрирования дифференциальных уравнений. Рассмотрим сначала процесс распространения погрешностей, внесенных в начальные значения. Пусть y_0^* — возмущенное начальное значение, $\varepsilon_0 = y_0 - y_0^*$ — его погрешность, а $y^*(t)$ — решение соответствующей задачи Коши

$$(y^*)'(t) = f(t, y^*(t)), \quad (14.7)$$

$$y^*(t_0) = y_0^*.$$

Вычтем из уравнения (14.1) уравнение (14.7) и воспользуемся формулой конечных приращений Лагранжа:

$$f(t, y(t)) - f(t, y^*(t)) = \lambda(t)(y(t) - y^*(t)), \quad \lambda(t) = f'_y(t, \tilde{y}(t)),$$

где $\tilde{y}(t)$ — некоторое промежуточное между $y(t)$ и $y^*(t)$ значение.

В результате получим, что погрешность $\varepsilon(t) = y(t) - y^*(t)$ удовлетворяет дифференциальному уравнению

$$\varepsilon'(t) = \lambda(t)\varepsilon(t) \quad (14.8)$$

и начальному условию

$$\varepsilon(t_0) = \varepsilon_0. \quad (14.9)$$

Решение задачи (14.8), (14.9) выражается формулой

$$\varepsilon(t) = \varepsilon_0 \exp \left\{ \int_{t_0}^t \lambda(\tau) d\tau \right\}.$$

Таким образом, величина

$$C(t) = \exp \left\{ \int_{t_0}^t \lambda(\tau) d\tau \right\} = \exp \left\{ \int_{t_0}^t f'_y(\tau, \tilde{y}(\tau)) d\tau \right\}$$

играет в задаче Коши роль коэффициента роста погрешности.

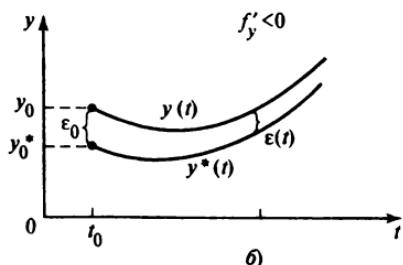
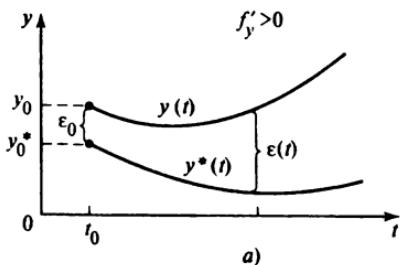


Рис. 14.3

Заметим, что знак производной f'_y оказывает существенное влияние на поведение погрешности $\varepsilon(t)$. Если $f'_y > 0$, то величина $C(t)$, а вместе с ней и модуль погрешности монотонно возрастают. При этом соответствующие интегральные кривые расходятся. Иллюстрацией такого поведения погрешности может служить рис. 14.3, а. Иначе ведет себя погрешность при $f'_y < 0$. Здесь $C(t)$ и $|\varepsilon(t)|$ с ростом t монотонно убывают, а соответствующие интегральные кривые сближаются. Погрешность, внесенная в начальное значение, имеет тенденцию к затуханию (рис. 14.3, б). Когда производная f'_y незнакопостоянна, поведение погрешности может быть более сложным.

Важно отметить, что в любом случае выполнение одностороннего условия Липшица (14.5) гарантирует, что коэффициент $C(t)$ роста погрешности окажется ограниченным, если задача решается на конечном отрезке $[t_0, T]$. В самом деле, в этой ситуации

$$\int_{t_0}^t f'_y(\tau, \tilde{y}(\tau)) d\tau \leq \int_{t_0}^t \sigma d\tau = \sigma(t - t_0)$$

и поэтому $C(t) \leq K(T)$ для всех $t_0 \leq t \leq T$, где

$$K(T) = \begin{cases} e^{\sigma(T-t_0)} & \text{при } \sigma > 0, \\ 1 & \text{при } \sigma \leq 0. \end{cases} \quad (14.10)$$

Таким образом, при выполнении условия $f'_y \leq \sigma$ справедлива оценка

$$\max_{t_0 \leq t \leq T} |y(t) - y^*(t)| \leq K(T) |y_0 - y_0^*|, \quad (14.11)$$

выражающая устойчивость на конечном отрезке $[t_0, T]$ решения задачи Коши по начальным значениям.

4. Модельное уравнение. Наиболее простым образом ведет себя погрешность, когда решается линейное уравнение

$$y'(t) = \lambda y(t) + f(t)$$

с постоянным коэффициентом λ . В этом случае погрешность ϵ удовлетворяет уравнению $\epsilon'(t) = \lambda \epsilon(t)$ и выражается формулой

$$\epsilon(t) = \epsilon_0 e^{\lambda(t-t_0)}. \quad (14.12)$$

Поскольку функция $f(t)$ не влияет на характер распространения погрешности, при изучении устойчивости по начальным значениям естественно ограничиться случаем $f(t) \equiv 0$ и рассматривать уравнение

$$y'(t) = \lambda y(t). \quad (14.13)$$

Уравнение (14.13) часто называют *модельным уравнением*. Оно играет важную роль при исследовании свойств численных методов решения задачи Коши.

Как следует из формулы (14.12), модуль погрешности решения уравнения (14.13) изменяется в e раз за интервал времени $\tau = \frac{1}{|\lambda|}$.

Поэтому величину $\tau = \frac{1}{|\lambda|}$ иногда называют *временной постоянной* или *постоянной времени* модельного уравнения (14.13). Если же параметр λ является комплексным числом, то временной постоянной называют величину $\tau = \frac{1}{|\operatorname{Re} \lambda|}$.

Когда рассматривается распространение малого возмущения, внесенного в решение $y(t)$ уравнения $y' = f(t, y)$ в малой окрестности точки \tilde{t} , значение коэффициента $\lambda(t)$ в уравнении (14.8) оказывается близко к постоянной $\tilde{\lambda} = f'_y(\tilde{t}, y(\tilde{t}))$. Поэтому при $t \approx \tilde{t}$ справедливо приближенное равенство $\epsilon'(t) \approx \tilde{\lambda} \epsilon(t)$. Это означает, что поведение погрешности для уравнения $y' = \tilde{\lambda} y$ моделирует локальное распространение погрешности для общего уравнения (14.1). Роль временной постоянной играет здесь *локальная временная постоянная* $\tau(\tilde{t}) = \frac{1}{|\tilde{\lambda}|} = \frac{1}{|f'_y(\tilde{t}, y(\tilde{t}))|}$,

значение которой меняется с изменением точки \tilde{t} . Для функции f , которая может принимать комплексные значения, формула для τ имеет вид

$$\tau(\tilde{t}) = \frac{1}{|\operatorname{Re} f'_y(\tilde{t}, y(\tilde{t}))|}.$$

5. Устойчивость по правой части. Будет ли решение задачи Коши устойчивым не только по отношению к погрешности ε_0 задания начального значения, но и к погрешностям $\psi(t)$ задания правой части уравнения? Положительный ответ на этот вопрос дает следующая теорема.

Теорема 14.3. Пусть выполнены условия теоремы 14.1. Далее, пусть $y(t)$ — решение задачи (14.1), (14.2), а $y^*(t)$ — решение задачи

$$(y^*)'(t) = f(t, y^*(t)) + \psi(t), \quad (14.14)$$

$$y^*(t_0) = y_0^*. \quad (14.15)$$

Тогда справедлива оценка

$$\max_{t_0 \leq t \leq T} |y(t) - y^*(t)| \leq K(T) \left(|y_0 - y_0^*| + \int_{t_0}^T |\psi(\tau)| d\tau \right), \quad (14.16)$$

выражающая устойчивость на конечном отрезке $[t_0, T]$ решения задачи Коши по начальным значениям и правой части. Здесь $K(T) = e^{L(T-t_0)}$. ■

Замечание. Величина $K(T)$ играет в задаче Коши роль оценки числа обусловленности¹. Если в теореме 14.1 условие Липшица (14.3) заменить односторонним условием (14.5), то оценка (14.16) будет выполнена с постоянной $K(T)$, определенной формулой (14.10).

6. Устойчивость решения на неограниченном промежутке. При решении самых разнообразных прикладных задач особый интерес представляет изучение описываемых дифференциальными уравнениями процессов на больших временных отрезках. В такой ситуации недостаточно наличие у задачи Коши свойства устойчивости на конечном отрезке. Если входящая в неравенство (14.16) величина $K(T)$ может неограниченно расти с ростом T , то это означает, что допускается неограниченный при $T \rightarrow \infty$ рост погрешностей. Как следствие, при достаточно больших T такая задача является плохо обусловленной и найти ее решение на отрезке $[t_0, T]$ с приемлемой точностью оказывается невозможно.

Пример 14.4. Рассмотрим задачу Коши $y'(t) = y(t) - \sin t + \cos t$, $y(0) = 0$. Ее решением, как нетрудно проверить, является функция $y(t) = \sin t$.

¹ Напомним, что общее понятие о числе обусловленности вычислительной задачи содержится в гл. 3.

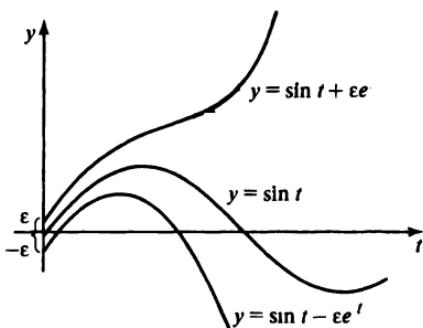


Рис. 14.4

Внесем в начальное значение погрешность, заменив условие $y(0) = 0$ условием $y(0) = \varepsilon$. Решением соответствующей задачи служит уже функция $y(t) = \sin t + \varepsilon e^t$. Погрешность εe^t с ростом t быстро увеличивается и, как видно из рис. 14.4, уже при не очень больших t ее значение становится неприемлемо большим. ▲

Для того чтобы обусловленность задачи Коши не ухудшалась с ростом T , в силу замечания 1 к теореме 14.3 достаточно потребовать, чтобы правая часть уравнения удовлетворяла неравенству $f'_y(t, y) \leq 0$ для всех $t \geq t_0$ и произвольных y . Более того, можно доказать, что при выполнении условия $f'_y(t, y) \leq \sigma < 0$ справедлива следующая оценка:

$$|y(t) - y^*(t)| \leq e^{\sigma(t-t_0)} |y_0 - y_0^*| + \frac{1}{|\sigma|} \max_{t_0 \leq t' \leq t} |\psi(t')|. \quad (14.17)$$

Предположим, что на каждом отрезке $[t_0, T]$ ($t_0 < T$ — произвольное) неравенство (14.6) выполнено с некоторой постоянной $\sigma = \sigma(T)$. Тогда решение $y(t)$ определено для всех $t_0 \leq t < \infty$. Пусть $y^*(t)$ — решение уравнения (14.7), отвечающее произвольному начальному значению y_0^* . Назовем решение задачи Коши (14.1), (14.2) *устойчивым по Ляпунову*¹, если справедлива оценка $\max_{t_0 \leq t \leq T} |y(t) - y^*(t)| \leq K |y_0 - y_0^*|$, где постоянная K не зависит от T . Если дополнительно известно, что $y^*(t) - y(t) \rightarrow 0$ при $t \rightarrow \infty$, то решение $y(t)$ называется *асимптотически устойчивым*.

¹ Александр Михайлович Ляпунов (1857—1918) — русский математик и механик. Приведенное здесь определение устойчивости по Ляпунову является более грубым, чем классическое определение [115].

Замечание 1. Решения модельного уравнения (14.13) с вещественным параметром λ устойчивы по Ляпунову тогда и только тогда, когда $\lambda \leq 0$, и асимптотически устойчивы тогда и только тогда, когда $\lambda < 0$. Этот вывод легко следует из формулы (14.12). Если же параметр λ — комплексное число, то из той же формулы следует, что $|\varepsilon(t)| = |\varepsilon_0| e^{\operatorname{Re}\lambda(T-t_0)}$. Поэтому решения модельного уравнения устойчивы по Ляпунову тогда и только тогда, когда $\operatorname{Re}\lambda \leq 0$, и асимптотически устойчивы тогда и только тогда, когда $\operatorname{Re}\lambda < 0$.

Замечание 2. Для решения задачи Коши (14.1), (14.2) (как вытекает из неравенства (14.11) и формулы (14.10)) грубым достаточным условием устойчивости по Ляпунову служит выполнение неравенства $f'_y \leq 0$. Следствием выполнения условия $f'_y \leq \sigma$ с постоянной $\sigma < 0$ является асимптотическая устойчивость решения.

§ 14.2. Численные методы решения задачи Коши. Основные понятия и определения

1. Сетки и сеточные функции. Первый этап на пути построения численного метода решения задачи Коши состоит в замене отрезка $[t_0, T]$ — области непрерывного изменения аргумента t — множеством $\bar{\omega}^h$, которое состоит из конечного числа точек $t_0 < t_1 < \dots < t_N = T$ и называется *сеткой*. Сами точки t_n называются *узлами сетки*, а величина $h_n = t_n - t_{n-1}$ — *шагом сетки* (рис. 14.5). Для того чтобы упростить изложение, будем рассматривать, как правило, *равномерные сетки*, т.е. такие сетки, для которых шаг h_n постоянен. В этом случае $h_n = h = (T - t_0)/N$ и $t_n = t_0 + nh$, $n = 1, 2, \dots, N$.

Наряду с функциями непрерывного аргумента будем рассматривать и *сеточные функции*, т.е. такие функции, которые определены лишь в узлах сетки $\bar{\omega}^h$. Для того чтобы отличать сеточные функции от функций непрерывного аргумента, будем помечать их индексом h . Так, например u^h — сеточная функция. Для краткости записи значения $u^h(t_n)$ сеточной функции u^h в узлах t_n сетки $\bar{\omega}^h$ будем обозначать через u_n .

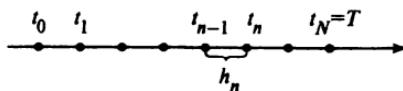


Рис. 14.5

2. Дискретная задача Коши. Следующий этап в построении численного метода состоит в замене задачи Коши ее дискретным аналогом — системой уравнений, решая которую можно последовательно находить значения y_1, y_2, \dots, y_N сеточной функции y^h , играющие роль приближений к значениям решения задачи Коши в узлах сетки $\bar{\omega}^h$.

В основе построения конкретного численного метода лежит тот или иной способ замены дифференциального уравнения $y' = f(t, y)$ его дискретным аналогом — уравнением вида

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y_{n+1-j} = \Phi(t_n, y_{n+1-k}, \dots, y_n, y_{n+1}, h), \quad (14.18)$$

в которое входят значения сеточной функции y^h в $k + 1$ последовательных точках $t_{n+1-k}, \dots, t_n, t_{n+1}$. Предполагается, что $\alpha_0 \neq 0$.

Во всех рассматриваемых в этой главе методах сумму

$$\frac{1}{h} (\alpha_0 y_{n+1} + \alpha_1 y_n + \dots + \alpha_k y_{n+1-k}), \quad (14.19)$$

стоящую в левой части уравнения (14.18), можно рассматривать как разностную аппроксимацию производной y' в соответствии с одной из формул численного дифференцирования (см. гл. 12). Правую часть Φ уравнения (14.18) можно рассматривать как специальным образом построенную аппроксимацию функции f .

Значение y_{n+1} приближенного решения в очередной точке находится из уравнения (14.18). При этом используются найденные ранее значения сеточной функции y^h в k предыдущих точках t_{n+1-k}, \dots, t_n . Поэтому такие методы получили название k -шаговых. Как нетрудно видеть, для того чтобы найти значения сеточной функции y^h во всех узлах сетки $\bar{\omega}^h$, используя k -шаговый метод, необходимо задать k начальных значений:

$$y^h(t_0) = y_0, \quad y^h(t_1) = y_1, \dots, y^h(t_{k-1}) = y_{k-1}. \quad (14.20)$$

Задачу вычисления сеточной функции y^h удовлетворяющей уравнению (14.18) для всех $n \geq k - 1$ и принимающей заданные начальные значения (14.20), будем называть *дискретной задачей Коши*.

Замечание. Принято считать, что уравнением (14.18) задается численный метод решения задачи Коши. Далее мы будем отождествлять свойства численного метода, дискретного уравнения (14.18) и соответствующей дискретной задачи Коши.

При $k = 1$ уравнение (14.18) упрощается и принимает вид

$$\frac{y_{n+1} - y_n}{h} = \Phi(t_n, y_n, y_{n+1}, h). \quad (14.21)$$

Соответствующий метод принято называть *одношаговыми*. Вычисление значения y_{n+1} осуществляется здесь с использованием только одного предыдущего значения y_n . Поэтому одношаговые методы часто называют *самостартующими*.

Пример 14.5. Простейший дискретный аналог дифференциального уравнения (14.1) представляет собой уравнение

$$\frac{y_{n+1} - y_n}{h} = f(t_n, y_n), \quad (14.22)$$

приводящее к известному *методу Эйлера*¹. ▲

Пример 14.6. Метод Эйлера является примером одношагового метода. Вычисление очередного значения y_{n+1} осуществляется здесь по формуле

$$y_{n+1} = y_n + hf(t_n, y_n). \quad (14.23)$$

При $k > 1$ численный метод называют *многошаговым*. Примеры таких методов можно найти в § 14.7 и 14.10.

Замечание. Использование многошагового метода предполагает преодоление одной специфической трудности, не возникающей при применении одношаговых методов. Как уже отмечалось выше, k -шаговый метод требует задания k начальных значений (14.20), в то время как в постановке задачи Коши содержится только одно начальное значение y_0 . Поэтому при $k > 1$ метод не является самостартующим, и для вычисления дополнительных значений y_1, y_2, \dots, y_{k-1} необходимы специальные подходы.

3. Явные и неявные методы. Реализация численного метода на компьютере предполагает построение алгоритма, позволяющего вычислить решение поставленной дискретной задачи Коши. В случае, когда входящая в уравнение (14.18) функция Φ не зависит от y_{n+1} , вычисление значения y_{n+1} не вызывает затруднений и осуществляется по явной формуле

$$y_{n+1} = \alpha_0^{-1} \left[-\sum_{j=1}^k \alpha_j y_{n+1-j} + h\Phi(t, y_{n+1-k}, \dots, y_n, h) \right],$$

¹ Леонард Эйлер (1707—1783) — математик, физик, механик, астроном. Родился в Швейцарии, с 1726 по 1741 г. и с 1776 по 1783 г. работал в России.

поэтому соответствующие методы называют **явными**. В противоположность им, методы, в которых функция Φ зависит от y_{n+1} , называют **неявными**. При реализации неявного метода при каждом n (или, как говорят, на каждом шаге) возникает необходимость решения относительно y_{n+1} нелинейного уравнения (14.18).

Пример 14.7. Метод Эйлера, для которого вычисления y_{n+1} производятся по явной формуле (14.23), представляет собой явный метод. ▲

Пример 14.8. Простейшим примером неявного метода является **неявный метод Эйлера**, соответствующий аппроксимации дифференциального уравнения (14.1) дискретным уравнением

$$\frac{y_{n+1} - y_n}{h} = f(t_{n+1}, y_{n+1}). \quad (14.24)$$

Другим примером неявного метода может служить *правило трапеций*

$$\frac{y_{n+1} - y_n}{h} = \frac{1}{2} (f(t_n, y_n) + f(t_{n+1}, y_{n+1})). \quad (14.25)$$

Как в том, так и в другом методе значение y_{n+1} определяется уравнением неявно, и для его вычисления приходится использовать один из итерационных методов решения нелинейных уравнений. ▲

4. Устойчивость. Если решение дискретной задачи Коши не обладает устойчивостью по отношению к малым возмущениям начальных значений и правой части уравнения, то соответствующий численный метод нельзя использовать в практических вычислениях. Приведем определение устойчивости, достаточное для понимания основного содержания этого и следующих четырех параграфов. Более подробно обсуждение этой проблемы будет проведено в § 14.8.

Внесем в правую часть уравнения (14.18) и в начальные условия (14.20) произвольные малые возмущения ψ_n и $\epsilon_0, \epsilon_1, \epsilon_2, \dots, \epsilon_{k-1}$ соответственно. Положим $y_0^* = y_0 - \epsilon_0, y_1^* = y_1 - \epsilon_1, \dots, y_{k-1}^* = y_{k-1} - \epsilon_{k-1}$. Пусть y^{*h} — решение соответствующей возмущенной задачи

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y_{n+1-j}^* = \Phi(t_n, y_n^*, y_{n+1-k}^*, \dots, y_n^*, y_{n+1}^*, h) + \psi_n, \quad (14.26)$$

$$y^{*h}(t_0) = y_0^*, y^{*h}(t_1) = y_1^*, \dots, y^{*h}(t_{k-1}) = y_{k-1}^*. \quad (14.27)$$

Будем называть дискретную задачу Коши (14.18), (14.20) и соответствующий численный метод *устойчивыми на конечном отрезке* (или просто *устойчивыми*), если при всех $h \leq h_0$ (где h_0 достаточно мало) справедливо неравенство

$$\max_{0 \leq n \leq N} |y_n - y_n^*| \leq \bar{K}(T) \left[\max_{0 \leq n \leq k-1} |\varepsilon_n| + \sum_{n=k-1}^{N-1} |\psi_n| \cdot h \right], \quad (14.28)$$

где величина \bar{K} не зависит от h , ε_n ($0 \leq n \leq k-1$) и ψ_n ($k-1 \leq n \leq N-1$).

Замечание. Неравенство (14.28) является дискретным аналогом неравенства (14.16), выражающего устойчивость решения задачи Коши. Для одношаговых методов (т.е. при $k = 1$) неравенство (14.28) принимает вид

$$\max_{0 \leq n \leq N} |y_n - y_n^*| \leq \bar{K}(T) \left(|\varepsilon_0| + \sum_{n=0}^{N-1} |\psi_n| \cdot h \right),$$

и аналогия с (14.16) становится еще более очевидной. Действительно, сумму $\sum_{n=0}^{N-1} |\psi_n| \cdot h$ можно рассматривать как дискретный аналог интеграла $\int_{t_0}^T |\psi(t)| dt$, построенный по формуле левых прямоугольников (см. § 13.1).

5. Аппроксимация. Пусть $y(t)$ — произвольная гладкая функция. Зафиксируем значение $t = t_n$ и устремим h к нулю (а n соответственно — к бесконечности). Будем предполагать, что замена в формуле (14.19) значений y_{n+1-j} сеточной функции y^h соответствующими значениями $y(t - (j-1)h)$ функции y дает величину

$$\frac{1}{h} (\alpha_0 y(t+h) + \alpha_1 y(t) + \dots + \alpha_k y(t-(k-1)h)), \quad (14.29)$$

стремящуюся к $y'(t)$ при $h \rightarrow 0$.

Аналогично предположим, что

$$\Phi(t, y(t-(k-1)h), \dots, y(t), y(t+h), h) \rightarrow f(t, y(t)) \text{ при } h \rightarrow 0.$$

Замечание. Из сделанных предположений следует, что коэффициенты $\alpha_0, \alpha_1, \dots, \alpha_k$ должны удовлетворять условию

$$\alpha_0 + \alpha_1 + \dots + \alpha_k = 0. \quad (14.30)$$

Действительно, для $y(t) \equiv 1$ величина (14.29) превращается в $\bar{\alpha}/h$, где $\bar{\alpha} = \sum_{j=0}^k \alpha_j$. По условию, $\bar{\alpha}/h \rightarrow y'(t) = 0$ при $h \rightarrow 0$. Но это возможно лишь при $\bar{\alpha} = 0$, что эквивалентно равенству (14.30).

Пусть $y(t)$ — решение задачи Коши (14.1), (14.2). Назовем сеточную функцию ψ^h , определяемую формулой

$$\psi_n = \frac{1}{h} \sum_{j=0}^k \alpha_j y(t_{n+1-j}) - \Phi(t_n, y(t_{n+1-k}), \dots, y(t_{n+1}), h),$$

погрешностью аппроксимации дискретного уравнения (14.18) на решении y . Эта формула, записанная в виде

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y(t_{n+1-j}) = \Phi(t_n, y(t_{n+1-k}), \dots, y(t_{n+1}), h) + \psi_n, \quad (14.31)$$

позволяет заметить, что функция $y(t)$ удовлетворяет уравнению (14.18) с точностью до погрешности аппроксимации ψ_n .

Сеточную функцию ψ^h используют для предварительной оценки того, насколько точно аппроксимируется дифференциальное уравнение его дискретным аналогом. Говорят, что *дискретное уравнение* (14.18) *аппроксимирует дифференциальное уравнение* (14.1), если $\max_{k-1 \leq n < N} |\psi_n| \rightarrow 0$ при $h \rightarrow 0$, и *аппроксимирует его с p-м порядком*,

если справедлива оценка $\max_{k-1 \leq n < N} |\psi_n| \leq Ch^p$, $p > 0$.

Часто для оценки качества одношаговых методов (14.21) используют не погрешность аппроксимации, а другую величину — локальную погрешность. Пусть y_{n+1} — значение, найденное из уравнения

$$\frac{y_{n+1} - y(t_n)}{h} = \Phi(t_n, y(t_n), y_{n+1}, h), \quad (14.32)$$

т.е. из уравнения (14.21), в которое вместо y_n подставлено точное значение решения дифференциального уравнения в точке $t = t_n$. Тогда разность $I_n = y(t_{n+1}) - y_{n+1}$ называется *локальной погрешностью метода* (или *его погрешностью на шаге*). Другими словами, I_n — это погрешность, которую допускает за один шаг метод, стартовавший с точного решения.

Когда Φ не зависит от y_{n+1} (т.е. метод (14.21) является явным), локальная погрешность и погрешность аппроксимации оказываются связаны простым равенством $I_n = \psi_n h$, что непосредственно вытекает из данных определений.

Пример 14.9. Покажем, что метод Эйлера имеет первый порядок аппроксимации.

Известно, что

$$\frac{y(t_n + h) - y(t_n)}{h} = y'(t_n) + \frac{h}{2} y''(\xi_n),$$

где $t_n < \xi_n < t_{n+1}$ (см. § 12.1). Учитывая равенство $y'(t_n) = f(t_n, y(t_n))$, для погрешности аппроксимации получаем следующее выражение:

$$\psi_n = \frac{y(t_n + h) - y(t_n)}{h} - f(t_n, y(t_n)) = \frac{h}{2} y''(\xi_n). \quad (14.33)$$

Поэтому

$$\max_{0 \leq n \leq N} |\psi_n| \leq \frac{M_2}{2} h, \quad M_2 = \max_{[t_0, T]} |y''(t)|, \quad (14.34)$$

т.е. метод действительно имеет первый порядок аппроксимации. ▲

Пример 14.10. Для погрешности аппроксимации неявного метода Эйлера (14.24) также справедлива оценка (14.34) и поэтому он также имеет первый порядок аппроксимации. Простое доказательство этого факта рекомендуем провести в качестве упражнения. ▲

Пример 14.11. Найдем выражение для локальной погрешности метода Эйлера.

По определению, $I_n = y(t_{n+1}) - y_{n+1}$, где $y_{n+1} = y(t_n) + hf(t_n, y(t_n))$. Но в силу равенства (14.33) $y(t_{n+1}) = y(t_n) + hf(t_n, y(t_n)) + h\psi_n = y_{n+1} + h\psi_n$. Поэтому $I_n = h\psi_n = \frac{h^2}{2} y''(\xi_n)$. Таким образом, локальная погрешность метода Эйлера имеет второй порядок малости относительно шага h . ▲

6. Сходимость. Пусть $y(t)$ — решение задачи Коши. Назовем *глобальной погрешностью* (или просто *погрешностью*) численного метода сеточную функцию ε^h со значениями $\varepsilon_n = y(t_n) - y_n$ в узлах t_n . В качестве меры абсолютной погрешности метода примем величину $E(h) = \max_{0 \leq n \leq N} |y(t_n) - y_n|$.

Численный метод решения задачи Коши называют *сходящимся*, если для него $E(h) \rightarrow 0$ при $h \rightarrow 0$. Принято говорить, что *метод сходится с p-м порядком точности* (или *имеет p-й порядок точности*), если для погрешности справедлива оценка $E(h) \leq Ch^p$, $p > 0$.

Покажем теперь, что для устойчивого численного метода из наличия аппроксимации с порядком p следует сходимость с тем же порядком. Будем предполагать, что начальные значения y_1, y_2, \dots, y_{k-1} заданы с p -м порядком точности¹, т.е. верна оценка $\max_{1 \leq n \leq k-1} |y(t_n) - y_n| \leq C_0 h^p$.

Справедлива следующая основная теорема.

Теорема 14.4. Пусть численный метод устойчив на конечном отрезке и имеет порядок аппроксимации, равный p . Тогда если начальные значения y_1, \dots, y_{k-1} заданы с p -м порядком точности, то и метод сходится с p -м порядком точности. ■

Доказательство. Пусть ψ_n — погрешность аппроксимации. Положим $y_n^* = y(t_n)$. Равенство (14.31) позволяет утверждать, что сеточная функция y^{*h} является решением дискретной задачи Коши (14.26), (14.27). Устойчивость метода означает выполнение неравенства (14.28), которое в силу равенств $y_n^* = y(t_n)$ и $y(t_0) = y_0$ можно переписать так:

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq \bar{K}(T) \left[\max_{1 \leq n \leq k-1} |y(t_n) - y_n| + \sum_{n=k-1}^{N-1} |\psi_n| \cdot h \right]. \quad (14.35)$$

Учитывая, что

$$\sum_{n=k-1}^{N-1} |\psi_n| \cdot h \leq (T - t_0) \max_{0 \leq n \leq N} |\psi_n| \leq (T - t_0) Ch^p,$$

правую часть неравенства (14.35) оцениваем величиной $\bar{K}(T) (C_0 + (T - t_0)C) h^p = \bar{C} h^p$. Итак, $\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq \bar{C} h^p$. ■

6. Связь с задачей вычисления интеграла. Существует тесная связь между проблемой решения задачи Коши и задачей вычисления интеграла с переменным верхним пределом

$$y(t) = \int_{t_0}^t f(\tau) d\tau, \quad t_0 \leq t \leq T, \quad (14.36)$$

Действительно, вычисление интеграла (14.36) эквивалентно решению задачи Коши

$$y'(t) = f(t), \quad y(t_0) = 0, \quad (14.37)$$

являющейся частным случаем более общей задачи (14.1), (14.2).

¹ Для одношаговых методов это предположение излишне.

Таким образом, всякий численный метод решения задачи Коши порождает соответствующий метод численного интегрирования. Например, метод Эйлера $y_{n+1} = y_n + hf(t_n)$ приводит к формуле левых прямоугольников:

$$y(t_n) \approx y_n = h \sum_{i=0}^{n-1} f(t_i). \quad (14.38)$$

Неявный метод Эйлера $y_{n+1} = y_n + hf(t_{n+1})$ дает формулу правых прямоугольников:

$$y(t_n) \approx y_n = h \sum_{i=1}^n f(t_i),$$

а правило трапеций (14.25) приводит к известной формуле трапеций

$$y(t_n) \approx \frac{h}{2} \sum_{i=1}^n (f(t_{i-1}) + f(t_i)).$$

На примере формулы (14.38) легко увидеть различие между локальной и глобальной погрешностями. Локальная погрешность — это погрешность, допускаемая на одном элементарном отрезке, т.е.

$$l_n = \int_{t_n}^{t_{n+1}} f(t) dt - hf(t_n),$$

а глобальная погрешность — это результирующая погрешность, т.е.

$$\varepsilon_n = \int_{t_0}^{t_n} f(t) dt - h \sum_{i=0}^{n-1} f(t_i).$$

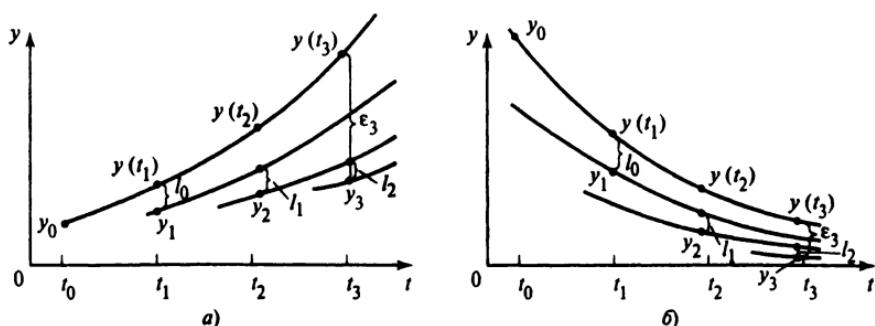


Рис. 14.6

В данном случае в силу линейности задачи (14.36) глобальная погрешность есть просто сумма локальных погрешностей: $\epsilon_n = \sum_{i=0}^{n-1} l_i$.

Для нелинейного уравнения $y' = f(t, y)$ это уже не так. В зависимости от характера поведения интегральных кривых глобальная погрешность может оказаться больше (рис. 14.6, а) или меньше (рис. 14.6, б) суммы соответствующих локальных погрешностей.

§ 14.3. Использование формулы Тейлора

Один из наиболее простых для понимания подходов к решению задачи Коши основан на использовании формулы Тейлора

$$y(t+h) = y(t) + y'(t)h + \frac{y''(t)}{2} h^2 + \dots + \frac{y^{(p)}(t)}{p!} h^p + R_{p+1}(t, h). \quad (14.39)$$

Здесь $R_{p+1}(t, h) = \frac{y^{(p+1)}(\xi)}{(p+1)!} h^{p+1}$ — остаточный член формулы Тейлора, ξ — некоторая точка, принадлежащая отрезку $[t, t+h]$.

Отбрасывая остаточный член, получаем приближенное равенство

$$y(t+h) \approx y(t) + y'(t)h + \frac{y''(t)}{2} h^2 + \dots + \frac{y^{(p)}(t)}{p!} h^p. \quad (14.40)$$

Если значение решения y в точке t известно, то в силу равенства

$$y'(t) = f(t, y(t)) \quad (14.41)$$

значение производной $y'(t)$ также можно считать известным. Для того чтобы вычислить производные y'', y''', \dots более высокого порядка, входящие в формулу (14.40), продифференцируем равенство (14.41) по t , используя правило дифференцирования сложной функции. Тогда получим

$$y'' = f'_t + f'_y y' = f'_t + f'_y f, \quad (14.42)$$

$$\begin{aligned} y''' &= f_{tt}^{(2)} + f_{ty}^{(2)} y' + (f_{yt}^{(2)} + f_{yy}^{(2)} y') f + f'_y (f'_t + f'_y y') = \\ &= f_{tt}^{(2)} + 2f_{ty}^{(2)} f + f'_y f'_t + (f'_y)^2 f + f_{yy}^{(2)} f^2 \end{aligned} \quad (14.43)$$

и т.д. Как нетрудно заметить, выражения для производных $y^{(k)}$ усложняются по мере роста порядка k .

Использование приближенной формулы (14.40) приводит к следующему явному одношаговому методу:

$$y_{n+1} = y_n + y'_n h + \frac{y''_n}{2!} h^2 + \dots + \frac{y^{(p)}_n}{p!} h^p. \quad (14.44)$$

Здесь $y'_n = f(t_n, y_n)$, значения y''_n и y'''_n получаются в результате подстановки в формулы (14.42) и (14.43) значений $t = t_n$ и $y = y_n$; аналогично вычисляются значения $y^{(k)}_n$ при $k > 3$.

Локальная погрешность этого метода I_n совпадает с величиной $R_{p+1}(t_n, h)$ остаточного члена формулы Тейлора, пропорциональной h^{p+1} . Воспользовавшись этим, можно доказать, что метод (14.44) сходится и имеет порядок точности, равный p .

Несмотря на то, что рассматриваемый метод теоретически дает возможность найти решение с любым порядком точности, на практике он применяется довольно редко. Дело в том, что использование формулы (14.44) приводит к необходимости вычисления большого числа частных производных $\frac{\partial^{l+s}f}{\partial t^l \partial y^s}$, что чаще всего является весьма трудоемкой и нередко аналитически невыполнимой операцией.

Более существенный аргумент против использования метода (14.44) состоит в том, что к настоящему времени разработаны эффективные численные методы решения задачи Коши (например, методы Рунге—Кутты и Адамса), предполагающие необходимость вычисления значений только функции f и не использующие ее частные производные. Именно этим методам, реализованным в виде стандартных программ и пакетов прикладных программ, мы и уделим основное внимание в дальнейшем.

Тем не менее для решения некоторых специальных классов задач приведенный выше метод может быть полезен. В частности, он используется при решении некоторых задач небесной механики, в которых вычисление производных $y^{(k)}$ не требует существенных дополнительных затрат в силу специальной структуры правых частей.

Пример 14.12. Найдем численное решение задачи Коши

$$y' = 2ty, \quad y(0) = 1 \quad (14.45)$$

на отрезке $[0, 1]$, использовав метод (14.44) при $p = 2$, обладающий вторым порядком точности. Как нетрудно проверить, точным решением этой задачи является функция $y(t) = e^{t^2}$.

Дифференцируя уравнение по t , получаем следующее выражение для второй производной: $y'' = 2y + 2ty' = 2(1 + 2t^2)y$, поэтому расчетная формула (14.44) в данном случае примет вид

$$y_{n+1} = y_n + h2t_n y_n + \frac{h^2}{2} 2(1 + 2t_n^2) y_n = (1 + 2t_n h + (1 + 2t_n^2)h^2) y_n.$$

Таблица 14.1

t	Метод (14.44) $h=0.1$	Точное решение	t	Метод (14.44) $h=0.1$	Точное решение
0.1	1.01000	1.01005	0.6	1.42840	1.43333
0.2	1.04050	1.04081	0.7	1.62438	1.63232
0.3	1.09336	1.09417	0.8	1.88396	1.89648
0.4	1.17186	1.17351	0.9	2.22835	2.24791
0.5	1.28108	1.28403	1.0	2.68783	2.71828

Найденное по этой формуле для $h = 0.1$ приближенное решение приведено в табл. 14.1 Для сравнения в ней же приведены значения точного решения.

Как видно из таблицы, решение оказалось найдено с точностью $\epsilon \approx 3 \cdot 10^{-2}$. ▲

§ 14.4. Метод Эйлера

1. **Геометрическая интерпретация метода Эйлера.** Простейшим и исторически первым численным методом решения задачи Коши (14.1), (14.2) является метод Эйлера. Его можно получить, если в приближенном равенстве (14.44) оставить только два первых слагаемых (т.е. взять $p = 1$). Тогда формула (14.44) примет вид

$$y_{n+1} = y_n + hf(t_n, y_n). \quad (14.46)$$

Геометрическая интерпретация одного шага метода Эйлера заключается в аппроксимации решения на отрезке $[t_n, t_{n+1}]$ касательной $y = y_n + + y'(t_n)(t - t_n)$, проведенной в точке (t_n, y_n) к интегральной кривой, проходящей через эту точку (рис. 14.7).

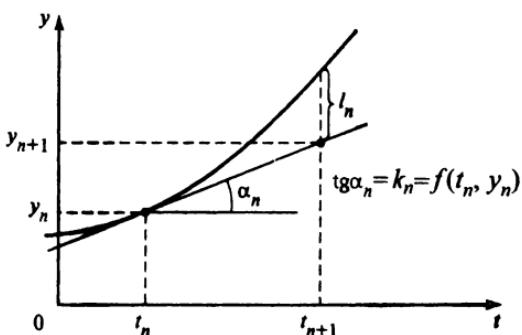


Рис. 14.7

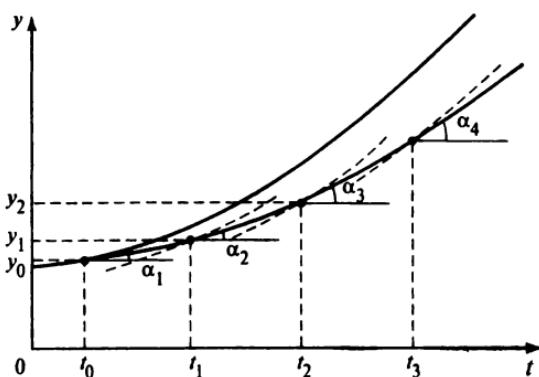


Рис. 14.8

Таким образом, после выполнения N шагов неизвестная интегральная кривая заменяется ломаной линией (ломаной Эйлера), для которой угловой коэффициент k_n очередного n -го звена равен значению $f(t_n, y_n)$ (рис. 14.8).

Как уже было отмечено в § 14.2, метод Эйлера представляет собой явный одношаговый метод. Для него погрешность аппроксимации имеет вид

$$\psi_n = \frac{h}{2} y''(\xi_n).$$

2. Устойчивость. Докажем, что метод Эйлера устойчив на конечном отрезке (см. определение устойчивости в § 14.2). Предварительно установим справедливость следующего вспомогательного утверждения.

Лемма 14.1. Пусть z^h — неотрицательная сеточная функция удовлетворяющая для всех $n \geq 0$ неравенству $z_{n+1} \leq (1 + \alpha)z_n + \beta_n$, где $\alpha \geq 0$, $\beta_n \geq 0$. Тогда при всех $n \geq 0$ справедлива оценка

$$z_n \leq e^{n\alpha} \left(z_0 + \sum_{k=0}^{n-1} \beta_k \right). \quad (14.47)$$

Доказательство. Справедливость неравенства (14.47) установим методом индукции. При $n = 0$ оно превращается в очевидное $z_0 \leq z_0$.

Пусть теперь неравенство (14.47) выполнено при некотором $n = m$. Тогда, используя оценки $1 + \alpha \leq e^\alpha$ и $1 \leq e^{(m+1)\alpha}$, получим следующую цепочку неравенств

$$z_{m+1} \leq (1 + \alpha)z_m + \beta_m \leq e^\alpha z_m + e^{(m+1)\alpha} \beta_m \leq$$

$$\leq e^{\alpha} e^{mh} \left(z_0 + \sum_{k=0}^{m-1} \beta_k \right) + e^{(m+1)\alpha} \beta_m = e^{(m+1)\alpha} \left(z_0 + \sum_{k=0}^m \beta_k \right),$$

т.е. неравенство (14.47) верно и при $n = m + 1$. Итак, оно верно при всех n . ■

Пусть теперь y^{*h} — решение возмущенной дискретной задачи Коши

$$y_{n+1}^* = y_n^* + h(f(t_n, y_n^*) + \psi_n), \quad (14.48)$$

$$y^{*h}(t_0) = y_0^*. \quad (14.49)$$

Теорема 14.5. Пусть функция f удовлетворяет условию $|f'_y| \leq L$. Тогда справедливо неравенство

$$\max_{0 \leq n \leq N} |y_n^* - y_n| \leq e^{L(T-t_0)} \left(|y_0^* - y_0| + h \sum_{k=0}^{N-1} |\psi_k| \right), \quad (14.50)$$

означающее, что метод Эйлера устойчив на конечном отрезке.

Доказательство. Вычитая из уравнения (14.48) уравнение (14.46) и пользуясь формулой конечных приращений Лагранжа

$$f(t_n, y_n^*) - f(t_n, y_n) = f'_y(t_n, \tilde{y}_n)(y_n^* - y_n),$$

получаем равенство

$$y_{n+1}^* - y_{n+1} = (1 + hf'_y(t_n, \tilde{y}_n))(y_n^* - y_n) + h\psi_n,$$

откуда следует

$$|y_{n+1}^* - y_{n+1}| \leq (1 + hL) |y_n^* - y_n| + h |\psi_n|. \quad (14.51)$$

Заметим теперь что для $z_n = |y_n^* - y_n|$, $\alpha = hL$ и $\beta_n = h|\psi_n|$ в силу (14.51) справедливо неравенство

$$z_{n+1} \leq (1 + \alpha)z_n + \beta_n.$$

Согласно лемме 14.1, имеем

$$|y_n^* - y_n| \leq e^{nhL} \left(|y_0^* - y_0| + \sum_{k=0}^{n-1} |\psi_k| h \right).$$

Учитывая, что $nh \leq Nh = T - t_0$, приходим к неравенству (14.50). ■

3. Оценка погрешности. Так как метод Эйлера устойчив на конечном отрезке и имеет первый порядок аппроксимации (см. пример 14.9), то из теоремы 14.4 следует, что он сходится с первым порядком точности. Точнее, верна следующая теорема.

Теорема 14.6. Пусть функция f удовлетворяет условию $|f'_y| \leq L$. Тогда для метода Эйлера справедлива такая оценка глобальной погрешности:

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq C(T)h, \quad (14.52)$$

где $C(T) = e^{L(T-t_0)}(T-t_0)M_2/2$, $M_2 = \max_{[t_0, T]} |y''(t)|$.

Приведем доказательство теоремы, не использующее теорему 14.4.

Доказательство. Пусть ψ_n — погрешность аппроксимации. Перепишем ее определение (14.33) в виде

$$y(t_{n+1}) = y(t_n) + h(f(t_n, y(t_n)) + \psi_n).$$

Полагая $y_n^* = y(t_n)$, замечаем, что сеточная функция y^{*h} является решением дискретной задачи Коши (14.48), (14.49), где $y_0^* = y_0$. Тогда в силу теоремы 14.5 справедлива оценка

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq e^{L(T-t_0)} \sum_{n=0}^{N-1} |\psi_n| h.$$

Учитывая, что $\sum_{n=0}^{N-1} |\psi_n| h \leq \max_{0 \leq n < N} |\psi_n|(T-t_0)$, и используя оценку (14.34) для погрешности аппроксимации, получаем неравенство (14.52). ■

Пример 14.13. Найдем численное решение задачи Коши $y' = 2ty$, $y(0) = 1$ на отрезке $[0, 1]$, использовав метод Эйлера. Заметим, что та же задача другим методом была решена в примере 14.12.

В данном случае расчетная формула (14.46) принимает вид

$$y_{n+1} = y_n + h2t_n y_n = (1 + 2ht_n)y_n. \quad (14.53)$$

Полученные с помощью этой формулы для значений шагов $h = 0.1$, $h = 0.01$ и $h = 0.001$ приближенные решения приведены в табл. 14.2. Для сравнения в последнем столбце даны значения точного решения $y(t) = e^{t^2}$. Нижняя строка таблицы содержит значения абсолютной погрешности $E(h)$. Как и следовало ожидать, при уменьшении шага в 10 раз погрешность уменьшается также примерно в 10 раз. ▲

Таблица 14.2

t	Метод Эйлера			Точное решение
	$h = 0.1$	$h = 0.01$	$h = 0.001$	
0.1	1.00000	1.00903	1.00995	1.01005
0.2	1.02000	1.03868	1.04060	1.04081
0.3	1.06080	1.09071	1.09383	1.09417
0.4	1.12445	1.16835	1.17299	1.17351
0.5	1.21440	1.27659	1.28328	1.28403
0.6	1.33584	1.42277	1.43226	1.43333
0.7	1.49615	1.61733	1.63080	1.63232
0.8	1.70561	1.87513	1.89432	1.89648
0.9	1.97850	2.21724	2.24480	2.24791
1.0	2.33463	2.67379	2.71376	2.71828
$E(h)$	$3.9 \cdot 10^{-1}$	$4.5 \cdot 10^{-2}$	$4.6 \cdot 10^{-3}$	—

4. Влияние вычислительной погрешности. Оценивая метод Эйлера необходимо учитывать, что при его реализации на компьютере неизбежно возникнут погрешности округления. В результате фактически вычисляемые значения y_n^* будут удовлетворять соотношению

$$y_{n+1}^* = y_n^* + h f(t_n, y_n^*) + \delta_n.$$

Величины δ_n учитывают вклад погрешностей округления. Это соотношение можно рассматривать как возмущенное уравнение вида (14.48), в котором $\psi_n = \delta_n/h$. Тогда неравенство (14.50) дает следующую оценку влияния погрешностей округления:

$$\max_{0 \leq n \leq N} |y_n^* - y_n| \leq e^{L(T-t_0)} \sum_{n=0}^{N-1} |\delta_n| \approx e^{L(T-t_0)} N \delta. \quad (14.54)$$

Здесь δ — некоторое среднее значение величины $|\delta_n|$, а $N = (T - t_0)/h$.

Таким образом, с учетом неравенств (14.52) и (14.54) получается следующая оценка погрешности фактически вычисляемых значений y_n^* :

$$\max_{0 \leq n \leq N} |y(t_n) - y_n^*| \lesssim \bar{E}_T(h) = e^{L(T-t_0)} (T-t_0) \left[\frac{M_2}{2} h + \frac{\delta}{h} \right].$$

Схематически график функции $\bar{E}_T(h)$ приведен на рис. 14.9.

Оказывается, что полная погрешность убывает только лишь при уменьшении шага h до некоторого значения $h_{\text{опт}}$. Достигимая точность

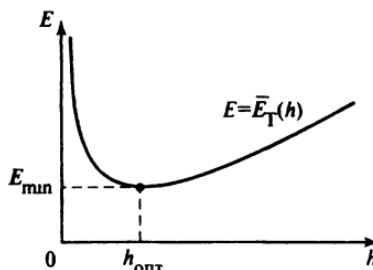


Рис. 14.9

метода ограничена снизу значением E_{\min} и попытка увеличить точность за счет уменьшения шага h при $h < h_{\text{опт}}$ приводит лишь к резкому росту погрешности.

Значение $h_{\text{опт}}$, как правило, бывает очень трудно определить заранее. Однако если очень высокая точность не нужна, то необходимый для ее достижения шаг h обычно бывает много больше $h_{\text{опт}}$.

Если же требуется высокая точность решения, то достичнуть ее с помощью метода Эйлера нельзя, даже если пойти на значительные затраты машинного времени (неизбежные при расчете с малым значением шага h). Необходимо иметь в своем распоряжении методы, имеющие более высокий порядок точности и позволяющие вести расчет со сравнительно крупным шагом h .

§ 14.5. Модификации метода Эйлера второго порядка точности

Медленная сходимость метода Эйлера (его погрешность убывает пропорционально лишь первой степени h) является серьезным препятствием для использования его на практике. Из рис. 14.7 видно, что уже один шаг по касательной к интегральной кривой приводит к значительной величине локальной погрешности I_n . Можно ли так подправить расчетную формулу метода, чтобы существенно уменьшить значение I_n ?

Пусть $y(t)$ — решение дифференциального уравнения

$$y'(t) = f(t, y(t)), \quad (14.55)$$

удовлетворяющее условию $y(t_n) = y_n$. Далее, пусть

$$k_n = \frac{y(t_{n+1}) - y(t_n)}{h} = \operatorname{tg} \gamma_n \quad (14.56)$$

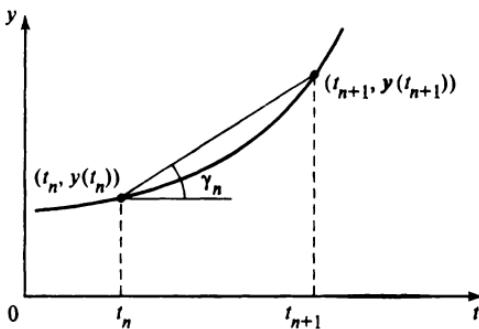


Рис. 14.10

— угловой коэффициент секущей, проходящей через точки \$(t_n, y(t_n))\$ и \$(t_{n+1}, y(t_{n+1}))\$ графика функции \$y(t)\$ (рис. 14.10).

Ясно, что «метод», состоящий в вычислении по формуле

$$y_{n+1} = y_n + h k_n, \quad (14.57)$$

имеет нулевую локальную погрешность. Для того чтобы воспользоваться этой формулой, нужно лишь «научиться вычислять значение \$k_n\$».

Интегрируя обе части уравнения (14.55) по \$t\$ от \$t_n\$ до \$t_{n+1}\$ и используя формулу Ньютона—Лейбница $\int_{t_n}^{t_{n+1}} y'(t) dt = y(t_{n+1}) - y(t_n)$, приходим к равенству

$$y(t_{n+1}) - y(t_n) = \int_{t_n}^{t_{n+1}} f(t, y(t)) dt. \quad (14.58)$$

Из равенств (14.56) и (14.58) следует, что

$$k_n = \frac{1}{h} \int_{t_n}^{t_{n+1}} f(t, y(t)) dt. \quad (14.59)$$

Заметим теперь, что применение для приближенного вычисления интеграла, стоящего в правой части выражения (14.59), формулы левых прямоугольников $\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx h f(t_n, y(t_n))$ немедленно приводит от (14.57) к методу Эйлера (14.46).

Известно (см. гл. 13), что больший порядок точности имеет формула трапеций

$$\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx \frac{h}{2} (f(t_n, y(t_n)) + f(t_{n+1}, y(t_{n+1}))).$$

Непосредственное ее применение к вычислению k_n приводит к правилу трапеций:

$$y_{n+1} = y_n + \frac{h}{2} (f(t_n, y_n) + f(t_{n+1}, y_{n+1})) \quad (14.60)$$

(ср. с (14.25)). Этот метод имеет второй порядок точности, но является неявным, поэтому его реализация связана с необходимостью решения относительно y_{n+1} нелинейного уравнения (14.60).

Построим на основе правила трапеций явный метод. Для этого представим в правую часть формулы (14.60) значение y_{n+1} , «предсказываемое» методом Эйлера. В результате получается метод

$$y_{n+1} = y_n + \frac{h}{2} (f(t_n, y_n) + f(t_{n+1}, y_n + hf(t_n, y_n))), \quad (14.61)$$

который называют *методом Эйлера—Коши* (или *методом Хьюна*).

Геометрическая иллюстрация этого метода представлена на рис 14.11. Вычисления разбиваются на два этапа. На первом этапе (этапе прогноза) в соответствии с методом Эйлера

$$y_{n+1}^{(0)} = y_n + hk_n^{(1)}, \quad k_n^{(1)} = f(t_n, y_n)$$

вычисляют грубое приближение к значению $y(t_{n+1})$. В точке $(t_{n+1}, y_{n+1}^{(0)})$ определяют угловой коэффициент $k_n^{(2)} = f(t_{n+1}, y_{n+1}^{(0)})$. На втором этапе

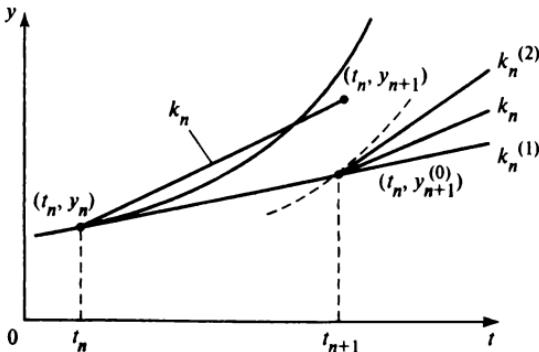


Рис. 14.11

(этапе коррекции) вычисляют усредненное значение углового коэффициента $k_n = (k_n^{(1)} + k_n^{(2)})/2$. Уточненное значение y_{n+1} находят по формуле $y_{n+1} = y_n + hk_n$, что соответствует шагу по прямой, проходящей через точку (t_n, y_n) и имеющей угловой коэффициент, равный k_n .

Замечание. Метод Эйлера—Коши относится к классу *методов прогноза и коррекции* (иначе говоря, *методов типа предиктор—корректор¹*).

Метод (14.61), который можно рассматривать как модификацию метода Эйлера, имеет второй порядок точности. Еще одну модификацию второго порядка точности можно получить с помощью формулы (центральных) прямоугольников

$$\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx hf(t_{n+1/2}, y(t_{n+1/2})), \quad t_{n+1/2} = t_n + h/2,$$

если для приближенного вычисления значения $y(t_{n+1/2})$ применить метод Эйлера. В результате получим расчетные формулы *усовершенствованного метода Эйлера*

$$\begin{aligned} y_{n+1/2} &= y_n + \frac{h}{2} k_n^{(1)}, \quad k_n^{(1)} = f(t_n, y_n); \\ y_{n+1} &= y_n + hk_n, \quad k_n = f(t_{n+1/2}, y_{n+1/2}). \end{aligned} \tag{14.62}$$

Геометрическая иллюстрация этого метода приведена на рис. 14.12.

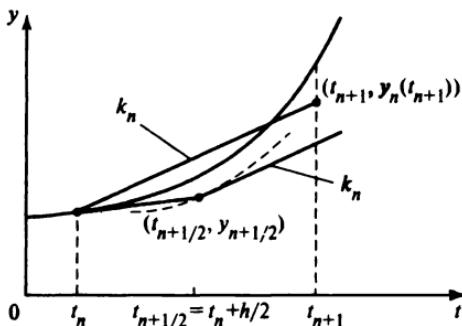


Рис. 14.12

¹ От английских слов to predict — «предсказывать», «прогнозировать»; to correct — «исправлять», «корректировать».

Пример 14.14. Применим рассмотренные в этом параграфе методы для численного решения задачи (14.45) с шагом $h = 0.1$.

Расчетная формула метода Эйлера—Коши принимает вид

$$\begin{aligned} y_{n+1} &= y_n + \frac{h}{2} (2t_n y_n + 2t_{n+1}(y_n + h t_n y_n)) = \\ &= y_n [1 + h(t_n + t_{n+1} + 2ht_n t_{n+1})]. \end{aligned} \quad (14.63)$$

Вычисления усовершенствованного метода Эйлера производим по формуле

$$y_{n+1} = y_n + 2h\left(t_n + \frac{h}{2}\right)\left(y_n + \frac{h}{2} 2t_n y_n\right) = y_n [1 + h(2t_n + h)(1 + ht_n)]. \quad (14.64)$$

Правило трапеций (14.60) приводит к уравнению

$$y_{n+1} = y_n + \frac{h}{2} (2t_n y_n + 2t_{n+1} y_{n+1}),$$

которое в данном случае линейно и легко разрешается относительно y_{n+1} :

$$y_{n+1} = y_n (1 + ht_n)(1 - ht_{n+1}). \quad (14.65)$$

Результаты вычислений по формулам (14.63)—(14.65) приведены в табл. 14.3. Там же для сравнения представлены значения решения $y(t) = e^{t^2}$. Нижняя строка таблицы содержит значения абсолютной погрешности $E(h)$.

Как видно из сравнения табл. 14.3 с табл. 14.2, проведенные в этом параграфе модификации метода Эйлера действительно привели к повышению точности. ▲

Таблица 14.3

t_n	Метод Эйлера—Коши; $h = 0.1$	Усовершенствованный метод Эйлера; $h = 0.1$	Правило трапеций; $h = 0.1$	Точное решение
0.1	1.01000	1.01000	1.01010	1.01005
0.2	1.04070	1.04060	1.04102	1.04081
0.3	1.09399	1.09367	1.09468	1.09417
0.4	1.17319	1.17253	1.17450	1.17351
0.5	1.28347	1.28228	1.28577	1.28403
0.6	1.43236	1.43038	1.43624	1.43333
0.7	0.63059	1.62749	1.63700	1.63232
0.8	1.89345	1.88870	1.90390	1.89648
0.9	2.24260	2.23546	2.25958	2.24791
1.0	2.70906	2.69843	2.73660	2.71828
$E(h)$	$0.93 \cdot 10^{-2}$	$2.0 \cdot 10^{-2}$	$1.9 \cdot 10^{-2}$	—

§ 14.6. Методы Рунге–Кутты

Наиболее популярными среди классических явных одношаговых методов являются методы Рунге–Кутты^{1,2}. Методы Эйлера, Эйлера–Коши и усовершенствованный метод Эйлера можно рассматривать как простейших представителей этого класса методов.

1. Вывод расчетных формул. Поясним вывод расчетных формул метода Рунге–Кутты. Пусть (как и в § 14.5) $y(t)$ — решение дифференциального уравнения $y' = f(t, y)$, удовлетворяющее условию $y(t_n) = y_n$. Запишем равенство (14.58) в следующем виде:

$$y(t_{n+1}) = y(t_n) + \int_{t_n}^{t_{n+1}} f(t, y(t)) dt. \quad (14.66)$$

Если бы входящий в это равенство интеграл можно было вычислить точно, то получилась бы простая формула, позволяющая последовательно вычислить значения решения в узлах сетки. Поскольку в действительности это невозможно, попробуем получить приближенную формулу, заменив интеграл квадратурной суммой (см. гл. 13).

На отрезке $[t_n, t_{n+1}]$ введем m вспомогательных узлов $t_n^{(1)} = t_n + \alpha_1 h$, $t_n^{(2)} = t_n + \alpha_2 h$, ..., $t_n^{(m)} = t_n + \alpha_m h$, где $0 = \alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_m \leq 1$. Заметим, что $t_n^{(1)} = t_n$, $t_n^{(m)} \leq t_{n+1}$. Заменяя входящий в равенство (14.66) интеграл квадратурной суммой с узлами $t_n^{(1)}, \dots, t_n^{(m)}$, получаем приближенное равенство

$$y(t_{n+1}) \approx y(t_n) + h \sum_{i=1}^m c_i f(t_n^{(i)}, y(t_n^{(i)})). \quad (14.67)$$

Однако воспользоваться равенством (14.67) для вычисления $y(t_{n+1})$ нельзя, так как значения функции y в точках $t_n^{(i)}$ для $i = 2, 3, \dots, m$ неизвестны. Чтобы найти эти значения, запишем равенства

$$y(t_n^{(i)}) = y(t_n) + \int_{t_n}^{t_n^{(i)}} f(t, y(t)) dt \quad (i=2, 3, \dots, m), \quad (14.68)$$

¹ Мартин Вильгельм Кутта (1867—1944) — немецкий математик.

² Первые методы этого типа Рунге предложил в 1895 г. Позже и независимо в 1901 г. Кутта дал общую схему вывода этих методов.

аналогичные равенства (14.66). Заменив для каждого i входящий в формулу (14.68) интеграл соответствующей ему квадратурной формулой с узлами $t_n^{(1)}, t_n^{(2)}, \dots, t_n^{(i-1)}$, придем к приближенным равенствам

$$y(t_n^{(2)}) \approx y(t_n) + h\beta_{21}f(t_n^{(1)}, y(t_n^{(1)})),$$

$$y(t_n^{(3)}) \approx y(t_n) + h(\beta_{31} f(t_n^{(1)}, y(t_n^{(1)})) + \beta_{32} f(t_n^{(2)}, y(t_n^{(2)}))),$$

$$y(t_n^{(i)}) \approx y(t_n) + h \sum_{j=1}^{i-1} \beta_{ij} f(t_n^{(j)}, y(t_n^{(j)})),$$

$$y(t_n^{(m)}) \approx y(t_n) + h \sum_{j=1}^{m-1} \beta_{mj} f(t_n^{(j)}, y(t_n^{(j)})),$$

позволяющим последовательно вычислять приближения к значениям $y(t_n^{(2)}), \dots, y(t_n^{(m)})$.

Обозначим теперь через $y_n^{(i)}$ вспомогательные величины, имеющие смысл приближений к значениям $y(t_n^{(i)})$; пусть $k_n^{(i)} = f(t_n^{(i)}, y_n^{(i)})$ — приближение к значению углового коэффициента k в точке $t_n^{(i)}$. Тогда расчетные формулы примут вид:

$$y_{n+1} = y_n + h k_n, \quad k_n = \sum_{i=1}^m c_i k_n^{(i)};$$

$$k_n^{(1)} = f(t_n^{(1)}, y_n^{(1)}), \quad y_n^{(1)} = y_n;$$

$$k_n^{(2)} = f(t_n^{(2)}, y_n^{(2)}), \quad y_n^{(2)} = y_n + h\beta_{21} k_n^{(1)};$$

$$k_n^{(3)} = f(t_n^{(3)}, y_n^{(3)}), \quad y_n^{(3)} = y_n + h(\beta_{31} k_n^{(1)} + \beta_{32} k_n^{(2)});$$

$$k_n^{(m)} = f(t_n^{(m)}, y_n^{(m)}), \quad y_n^{(m)} = y_n + h \sum_{j=1}^{m-1} \beta_{mj} k_n^{(j)}.$$

Часто из этих формул исключают вспомогательные величины $y_n^{(i)}$ и записывают формулы так:

$$y_{n+1} = y_n + h k_n, \quad k_n = \sum_{i=1}^m c_i k_n^{(i)};$$

$$k_n^{(1)} = f(t_n, y_n),$$

$$k_n^{(2)} = f(t_n + \alpha_2 h, y_n + h\beta_{21} k_n^{(1)}),$$

$$k_n^{(3)} = f(t_n + \alpha_3 h, y_n + h(\beta_{31} k_n^{(1)} + \beta_{32} k_n^{(2)})),$$

.....

$$k_n^{(m)} = f\left(t_n + \alpha_m h, y_n + h \sum_{j=1}^{m-1} \beta_{mj} k_n^{(j)}\right).$$

Заметим, что выведенные формулы задают явный одношаговый метод вида $y_{n+1} = y_n + h\Phi(t_n, y_n, h)$, где для вычисления значений функции $\Phi(t_n, y_n, h) = \sum_{i=1}^m c_i k_n^{(i)}$ используются значения правой части f в m вспомогательных точках. Поэтому этот метод называют **явным m -этапным методом Рунге—Кутты**.

Выбор конкретных значений параметров c_i , α_i , β_{ij} осуществляется, исходя из различных соображений. Естественно, что одним из основных является желание сделать порядок аппроксимации p максимально возможным.

2. Устойчивость и сходимость. Следующая теорема позволяет в дальнейшем называть методы Рунге—Кутты, имеющие p -й порядок аппроксимации, методами p -го порядка точности.

Теорема 14.7. Пусть правая часть дифференциального уравнения удовлетворяет условию $|f'_y| \leq L$. Тогда всякий явный m -этапный метод Рунге—Кутты устойчив на конечном отрезке. ■

Следствие. Пусть выполнено условие $|f'_y| \leq L$. Тогда если явный m -этапный метод Рунге—Кутты имеет p -й порядок аппроксимации, то он сходится с p -м порядком точности.

Справедливость следствия вытекает из теоремы 14.4.

3. Семейство явных двухэтапных методов. Выведем расчетные формулы семейства явных двухэтапных методов Рунге—Кутты второго порядка точности. Запишем формулы явного двухэтапного метода

$$y_{n+1} = y_n + h(c_1 k_n^{(1)} + c_2 k_n^{(2)}),$$

$$k_n^{(1)} = f(t_n, y_n), \quad k_n^{(2)} = f(t_n + \alpha h, y_n + h\beta k_n^{(1)})$$

в виде

$$\frac{y_{n+1} - y_n}{h} = c_1 f(t, y_n) + c_2 f(t_n + \alpha h, y_n + \beta h f(t_n, y_n)).$$

Параметрами этого метода являются величины c_1, c_2, α, β .

Представим погрешность аппроксимации

$$\psi = \frac{y(t+h) + y(t)}{h} - c_1 f(t, y) - c_2 f(t + \alpha h, y + h\beta f(t, y))$$

(где $t = t_n$, $y = y(t_n)$, $y(t_n)$ — решение дифференциального уравнения $y' = f(t, y)$) в виде разложения по степеням h .

Формула Тейлора

$$y(t+h) = y(t) + y'(t)h + \frac{y''(t)}{2} h^2 + O(h^3)$$

с учетом равенств $y' = f$, $y'' = f'_t + f'_y f$ (см. (14.42)) дает формулу

$$\frac{y(t+h) - y(t)}{h} = f + \frac{1}{2} (f'_t + f'_y f)h + O(h^2)$$

(аргументы t, y у функции f и ее частных производных f'_t, f'_y опускаем).

Представим значение функции f в точке $(t + \alpha h, y + h\beta f)$, используя формулу Тейлора для функции двух переменных с центром в точке (t, y) :

$$f(t + \alpha h, y + h\beta f) = f(t, y) + f'_t \alpha h + f'_y h\beta f + O(h^2).$$

Таким образом,

$$\psi = (1 - c_1 - c_2)f + [(1/2 - c_2\alpha)f'_t + (1/2 - c_2\beta)f'_y f]h + O(h^2).$$

Если потребовать, чтобы выполнялись условия $c_1 + c_2 = 1$, $c_2\alpha = 1/2$, $c_2\beta = 1/2$ (что эквивалентно выбору $c_1 = 1 - 1/(2\alpha)$, $c_2 = 1/(2\alpha)$, $\beta = \alpha$), то первые два слагаемых в формуле для ψ обращаются в нуль, и поэтому метод будет иметь второй порядок аппроксимации.

Итак (с учетом следствия из теоремы 14.7), можно утверждать, что при любом $\alpha \in (0, 1]$ метод

$$y_{n+1} = y_n + h \left[\left(1 - \frac{1}{2\alpha}\right) f(t_n, y_n) + \frac{1}{2\alpha} f(t_n + \alpha h, y_n + \alpha h f(t_n, y_n)) \right] \quad (14.69)$$

имеет второй порядок точности.

Заметим, что при $\alpha = 1$ формула (14.69) дает метод Эйлера—Коши, а при $\alpha = 1/2$ — усовершенствованный метод Эйлера (см. § 14.5).

4. Метод Рунге—Кутты четвертого порядка точности. Наиболее известным из методов Рунге—Кутты является классический четырехэтапный метод четвертого порядка точности:

$$\begin{aligned} y_{n+1} &= y_n + h k_n, \quad k_n = \frac{1}{6} (k_n^{(1)} + 2k_n^{(2)} + 2k_n^{(3)} + k_n^{(4)}); \\ k_n^{(1)} &= f(t_n, y_n), \quad k_n^{(2)} = f\left(t_n + \frac{h}{2}, y_n + \frac{h}{2} k_n^{(1)}\right), \\ k_n^{(3)} &= f\left(t_n + \frac{h}{2}, y_n + \frac{h}{2} k_n^{(2)}\right), \quad k_n^{(4)} = f(t_n + h, y_n + h k_n^{(3)}). \end{aligned} \quad (14.70)$$

Этот метод весьма прост и, как показывает практика, довольно эффективен в обычных расчетах, когда отрезок $[t_0, T]$ не очень велик и нужна сравнительно невысокая точность.

Замечание. Применение метода (14.70) к решению задачи о вычислении интеграла (14.36) порождает формулу Симпсона

$$y_{n+1} = y_n + \frac{h}{6} (f(t_n) + 4f(t_n + 1/2) + f(t_n + 1)).$$

Таким образом, классический метод Рунге—Кутты четвертого порядка точности (14.70) можно рассматривать как аналог формулы Симпсона, отвечающий решению задачи Коши.

Пример 14.15. Продемонстрируем работу метода Рунге—Кутты четвертого порядка точности применительно к решению задачи Коши (14.45).

В этом случае расчетные формулы принимают вид

$$\begin{aligned} y_{n+1} &= y_n + h k_n, \quad k_n = \frac{1}{6} (k_n^{(1)} + 2k_n^{(2)} + 2k_n^{(3)} + k_n^{(4)}); \\ k_n^{(1)} &= 2t_n y_n, \quad k_n^{(2)} = 2\left(t_n + \frac{h}{2}\right)\left(y_n + \frac{h}{2} k_n^{(1)}\right), \\ k_n^{(3)} &= 2\left(t_n + \frac{h}{2}\right)\left(y_n + \frac{h}{2} k_n^{(2)}\right), \quad k_n^{(4)} = 2(t_n + h)(y_n + h k_n^{(3)}). \end{aligned}$$

Найденные с шагом $h = 0.1$ приближенные значения решения y_n и их погрешности ε_n приведены в табл. 14.4. ▲

Таблица 14.4

t_n	y_n	ε_n	t_n	y_n	ε_n
0.1	1.010050167	10^{-9}	0.6	1.433328994	$5 \cdot 10^{-7}$
0.2	1.040810770	$4 \cdot 10^{-9}$	0.7	1.632315187	$2 \cdot 10^{-6}$
0.3	1.094174265	$2 \cdot 10^{-8}$	0.8	1.896478467	$3 \cdot 10^{-6}$
0.4	1.173510814	$6 \cdot 10^{-8}$	0.9	2.247902590	$6 \cdot 10^{-6}$
0.5	1.284025256	$2 \cdot 10^{-7}$	1.0	2.718270175	$2 \cdot 10^{-5}$

Есть и другие методы Рунге—Кутты четвертого порядка точности. В качестве примера приведем расчетные формулы одного из таких методов

$$y_{n+1} = y_n + hk_n, \quad k_n = \frac{1}{6} (k_n^{(1)} + 4k_n^{(3)} + k_n^{(4)}),$$

$$k_n^{(1)} = f(t_n, y_n), \quad k_n^{(2)} = f\left(t_n + \frac{1}{2}h, y_n + \frac{h}{2}k_n^{(1)}\right),$$

$$k_n^{(3)} = f\left(t_n + \frac{1}{2}h, y_n + h\left(\frac{1}{4}k_n^{(1)} + \frac{1}{4}k_n^{(2)}\right)\right),$$

$$k_n^{(4)} = f(t_n + h, y_n + h(-k_n^{(2)} + 2k_n^{(3)})).$$

5. Метод Рунге—Кутты—Фельберга. Для построения методов Рунге—Кутты пятого порядка точности необходимо уже $m = 6$ вспомогательных точек. Приведем один из таких методов, предложенный в 1970 г. Фельбергом

$$y_{n+1} = y_n + hk_n, \quad k_n = \sum_{i=1}^6 c_i k_n^{(i)},$$

$$c_1 = \frac{16}{135}, \quad c_2 = 0, \quad c_3 = \frac{6656}{12825}, \quad c_4 = \frac{28561}{56430}, \quad c_5 = -\frac{9}{50}, \quad c_6 = \frac{2}{55},$$

$$k_n^{(1)} = f(t_n, y_n),$$

$$k_n^{(2)} = f\left(t_n + \frac{1}{4}h, y_n + h\frac{1}{4}k_n^{(1)}\right),$$

$$k_n^{(3)} = f\left(t_n + \frac{3}{8}h, y_n + h\left(\frac{3}{32}k_n^{(1)} + \frac{9}{32}k_n^{(2)}\right)\right),$$

$$k_n^{(4)} = f\left(t_n + \frac{12}{13}h, y_n + h\left(\frac{1932}{2197}k_n^{(1)} - \frac{7200}{2197}k_n^{(2)} + \frac{7296}{2197}k_n^{(3)}\right)\right),$$

$$k_n^{(5)} = f\left(t_n + h, y_n + h\left(\frac{439}{216}k_n^{(1)} - 8k_n^{(2)} + \frac{3680}{513}k_n^{(3)} - \frac{845}{4104}k_n^{(5)}\right)\right),$$

$$k_n^{(6)} = f\left(t_n + \frac{1}{2}h, y_n + h\left(-\frac{8}{27}k_n^{(1)} + 2k_n^{(2)} - \frac{3544}{2565}k_n^{(3)} + \frac{1859}{4104}k_n^{(4)} - \frac{11}{40}k_n^{(5)}\right)\right).$$

Фельберг обнаружил также, что вычисленные значения угловых коэффициентов можно скомбинировать иначе, получив метод четвертого порядка точности

$$y_{n+1}^* = y_n + hk_n^*, \quad k_n^* = \sum_{i=1}^6 c_i^* k_n^{(i)},$$

$$c_1^* = \frac{25}{216}, \quad c_2^* = 0, \quad c_3^* = \frac{1408}{2565}, \quad c_4^* = \frac{2197}{4104}, \quad c_5^* = -\frac{1}{5}, \quad c_6^* = 0.$$

При этом разность $y_{n+1}^* - y_{n+1} = \sum_{i=1}^6 (c_i^* - c_i) k_n^{(i)}$ дает оценку локальной погрешности для метода четвертого порядка точности и находится практически без дополнительных вычислительных затрат. Таким образом, метод Рунге—Кутты—Фельберга дает не только приближенные значения решения с пятым порядком точности, но и оценки локальных погрешностей.

6. Метод Дормана—Принса. Несмотря на то, что метод Рунге—Кутты—Фельберга широко и успешно используется на практике, у него есть один существенный недостаток. Локальная погрешность, которую оценивает метод, относится не к основному методу, имеющему пятый порядок точности, а к вспомогательному методу четвертого порядка точности.

В 1980 году Дорманом и Принсом был построен метод пятого порядка точности, лишенный этого недостатка. Приведем расчетные формулы *метода Дормана—Принса*

$$y_{n+1} = y_n + h k_n, \quad k_n = \sum_{i=1}^6 c_i k_n^{(i)};$$

$$c_1 = \frac{35}{384}, \quad c_2 = 0, \quad c_3 = \frac{500}{1113}, \quad c_4 = \frac{125}{192}, \quad c_5 = -\frac{2187}{6784}, \quad c_6 = \frac{11}{84};$$

$$k_n^{(1)} = f(t_n, y_n),$$

$$k_n^{(2)} = f\left(t_n + \frac{1}{5}h, y_n + h \frac{1}{5}k_n^{(1)}\right),$$

$$k_n^{(3)} = f\left(t_n + \frac{3}{10}h, y_n + h \left(\frac{3}{40}k_n^{(1)} + \frac{9}{40}k_n^{(2)}\right)\right),$$

$$k_n^{(4)} = f\left(t_n + \frac{4}{5}h, y_n + h \left(\frac{44}{45}k_n^{(1)} - \frac{56}{15}k_n^{(2)} + \frac{32}{9}k_n^{(3)}\right)\right),$$

$$k_n^{(5)} = f\left(t_n + \frac{8}{9}h, y_n + h \left(\frac{19372}{6561}k_n^{(1)} - \frac{25360}{2187}k_n^{(2)} + \frac{64448}{6561}k_n^{(3)} - \frac{212}{729}k_n^{(5)}\right)\right),$$

$$k_n^{(6)} = f\left(t_n + h, y_n + h \left(\frac{9017}{3168}k_n^{(1)} - \frac{355}{33}k_n^{(2)} + \frac{46732}{5247}k_n^{(3)} + \frac{49}{176}k_n^{(4)} - \frac{5103}{18656}k_n^{(5)}\right)\right),$$

$$k_n^{(7)} = f(t_n + h, y_{n+1}).$$

Вычисленные значения угловых коэффициентов можно скомбинировать иначе, получив метод 4-го порядка точности

$$y_{n+1}^* = y_n + hk_n^*, \quad k_n^* = \sum_{i=1}^7 c_i^* k_n^{(i)}; \\ c_1^* = \frac{5179}{57600}, \quad c_2^* = 0, \quad c_3^* = \frac{7571}{16695}, \quad c_4^* = \frac{393}{640}, \quad c_5^* = -\frac{92097}{339200}, \\ c_6^* = \frac{187}{2100}, \quad c_7^* = \frac{1}{40}.$$

При этом разность $y_{n+1}^* - y_{n+1}$ дает оценку локальной погрешности метода.

Заметим, что $k_{n+1}^{(1)} = k_n^{(7)}$, поэтому на каждом шаге метода требуется лишь шесть новых вычислений значений функции f .

Таким образом, метод Дормана—Принса дает не только приближенные значения решения с пятым порядком точности, но и оценки локальных погрешностей этого метода. При решении прикладных задач этот метод оказался очень эффективным и постепенно вытесняет метод Рунге—Кутты—Фельберга из практики вычислений.

Методы Рунге—Кутты—Фельберга и Дормана—Принса являются представителями семейства *вложенных методов Рунге—Кутты*. В это семейство также входит весьма популярный метод Дормана—Принса восьмого порядка точности, в котором для контроля погрешности используется метод седьмого порядка точности. Дополнительную информацию о вложенных методах Рунге—Кутты можно найти, например, в [109].

7. Обсуждение методов Рунге—Кутты. Методы Рунге—Кутты имеют несколько достоинств, определивших их популярность среди значительного числа исследователей. Эти методы легко программируются¹. Они обладают достаточными для широкого круга задач свойствами точности и устойчивости. Эти методы (как и все одношаговые методы) являются самостартующими и позволяют на любом этапе вычислений легко изменять шаг интегрирования.

Увеличивая число m вспомогательных точек, можно построить методы Рунге—Кутты любого порядка точности p . Однако уже при $p > 5$ эти методы используются довольно редко. Это объясняется как чрез-

¹ Это очень важно, если отсутствуют или по каким-либо причинам недоступны соответствующие стандартные программы. Если же используются развитые пакеты прикладных программ, то сложность программирования метода не интересует пользователя, поскольку он обращает внимание на другие свойства метода.

мерной громоздкостью получающихся вычислительных формул, так и тем, что преимущества методов высокого порядка точности p над методами, в которых $p = 4$ и $p = 5$, проявляются либо в тех задачах, где нужна очень высокая точность и используются компьютеры высокой разрядности, либо в тех задачах, где решение очень гладкое. Кроме того, методы Рунге—Кутты высокого порядка точности p часто оказываются менее эффективными по сравнению с методами Адамса того же порядка точности (см. § 14.7).

Замечание. Кроме описанных выше классических явных методов Рунге—Кутты используются и более сложные в реализации *неявные m -этапные методы Рунге—Кутты*:

$$y_{n+1} = y_n + hk_n, \quad k_n = \sum_{i=1}^m c_i k_n^{(i)};$$

$$k_n^{(i)} = f\left(t_n + \alpha_i h, y_n + h \sum_{j=1}^m \beta_{ij} k_n^{(j)}\right), \quad i = 1, 2, \dots, m.$$

Эти методы имеют ряд преимуществ перед явными методами, однако это достигается за счет существенного усложнения вычислительного алгоритма, так как на каждом шаге необходимо решать систему m нелинейных уравнений. В настоящее время неявные методы Рунге—Кутты применяются в основном для решения так называемых жестких задач (см. § 14.11).

8. Автоматический выбор шага. Отметим, что в современных программах, реализующих методы Рунге—Кутты, обязательно используется некоторый алгоритм автоматического изменения шага интегрирования $h_{n+1} = t_{n+1} - t_n$.

Интуитивно ясно, что на участках плавного изменения решения счет можно вести с достаточно крупным шагом. В то же время на тех участках, где происходят резкие изменения поведения решения, необходимо выбирать мелкий шаг интегрирования. Обычно начальное значение шага h_1 задает пользователь. Далее шаг интегрирования меняется в соответствии со значением получаемой в ходе вычислений оценки локальной погрешности. Само по себе изменение шага h_n для методов Рунге—Кутты (впрочем, как и для всех других одношаговых методов) не представляет сложности. Действительная проблема состоит в том, как оценить локальную погрешность и выбрать очередной шаг интегрирования.

Один из распространенных подходов состоит в использовании *правила Рунге* (*правила двойного пересчета*). Пусть значение y_n в точке t_n

уже найдено и $y(t)$ — решение уравнения $y'(t) = f(t, y)$, удовлетворяющее условию $y(t_n) = y_n$. Обозначим через $y^{h_{n+1}}$ приближение y_{n+1} к значению $y(t_{n+1})$, найденное с помощью одношагового метода

$$y_{n+1} = y_n + h_{n+1} \Phi(t_n, y_n, y_{n+1}, h_{n+1}), \quad (14.71)$$

который имеет порядок точности, равный p . Можно показать, что для методов Рунге—Кутты локальная погрешность $I_n = y(t_{n+1}) - y^{h_{n+1}}$ допускает представление

$$I_n = r(t_n, y_n) h_{n+1}^{p+1} + o(h_{n+1}^{p+1}),$$

где $r(t, y)$ — непрерывная функция. Следовательно, при достаточно малых h_{n+1} справедливо приближенное равенство

$$y(t_{n+1}) - y^{h_{n+1}} \approx r(t_n, y_n) h_{n+1}^{p+1}. \quad (14.72)$$

Уменьшим теперь шаг интегрирования вдвое, положив $h_{n+1/2} = h_{n+1}/2$, и вычислим приближение к значению решения в точке t_{n+1} с помощью того же одношагового метода. Для этого потребуется выполнить уже два элементарных шага по формулам

$$\begin{aligned} y_{n+1/2} &= y_n + h_{n+1/2} \Phi(t_n, y_n, y_{n+1/2}, h_{n+1/2}), \\ y_{n+1} &= y_{n+1/2} + h_{n+1/2} \Phi(t_{n+1/2}, y_{n+1/2}, y_{n+1}, h_{n+1/2}). \end{aligned}$$

Полученное таким образом значение $y^{h_{n+1/2}} = y_{n+1}$ будет, конечно, отличаться от значения, найденного по формуле (14.71). Достаточно ясно, что два шага длины $h_{n+1/2}$ приведут здесь к локальной погрешности

$$y(t_{n+1}) - y^{h_{n+1/2}} \approx 2r(t_n, y_n) h_{n+1/2}^{p+1}. \quad (14.73)$$

Вычитая из равенства (14.72) равенство (14.73), получаем формулу

$$y^{h_{n+1/2}} - y^{h_{n+1}} \approx (2^p - 1) 2r(t_n, y_n) h_{n+1/2}^{p+1}.$$

Сравнение ее с (14.73) приводит к приближенному равенству

$$y(t_{n+1}) - y^{h_{n+1/2}} \approx \frac{y^{h_{n+1/2}} - y^{h_{n+1}}}{2^p - 1}. \quad (14.74)$$

Использование этой формулы для апостериорной оценки локальной погрешности значения $y^{h_{n+1/2}}$ (которое в дальнейшем принимается за приближенное значение решения задачи Коши в точке t_{n+1}) и назы-

вают *правилом Рунге*¹. Заметим, что этот способ контроля точности приводит к увеличению времени счета примерно на 50 %.

Существуют более экономичные методы оценки локальной погрешности, основанные на использовании для контроля точности двух различных методов Рунге—Кутты. В настоящее время одним из самых эффективных методов такого типа является метод Дормана—Принса (см. п. 6). В этом методе для оценки погрешности метода пятого порядка точности используются формулы метода четвертого порядка точности, причем на одном шаге требуется всего лишь шесть вычислений значений правой части².

После того как тем или иным способом оценена локальная погрешность, программа принимает решение о том, оставить ли шаг интегрирования прежним, уменьшить его в 2 раза или увеличить в 2 раза. Это происходит примерно по той же схеме, что и в адаптивных программах, предназначенных для вычисления определенных интегралов (см. § 13.5). Известно, что при оптимальном выборе шагов интегрирования абсолютные погрешности, приходящиеся на каждый из шагов, должны быть примерно равны (см. [9]). Этот результат учитывается при создании стандартных программ с автоматическим выбором шага.

9. Влияние вычислительной погрешности. Влияние погрешностей на результат вычислений с помощью явных методов Рунге—Кутты примерно таково же, как и для метода Эйлера (см. § 14.4). Однако для них $\bar{E}_T(h) = C(T) \left(M_{p+1} h^p + \frac{\delta}{h} \right)$. Кроме того, высокая точность методов позволяет вести интегрирование со сравнительно большим шагом $h \gg \gg h_{\text{опт}}$ и поэтому влияние вычислительной погрешности обычно бывает несущественным.

§ 14.7. Линейные многошаговые методы.

Методы Адамса

1. Методы Адамса. В одношаговых методах после того как найдено очередное значение y_n в точке t_n , значение y_{n-1} отбрасывают и уже не используют в последующих вычислениях. Естественно все же попытаться извлечь определенную пользу из информации о значениях решения $y_{n-k+1}, \dots, y_{n-1}, y_n$ не в одной, а в k предыдущих точках $t_{n-k+1}, \dots, t_{n-1}, t_n$, т.е. применить многошаговый метод.

¹ Отметим, что использование правила Рунге требует определенной осторожности, так как равенство (14.72) имеет место, вообще говоря, лишь при достаточно малом значении шага h_{n+1} .

² Более подробное описание этого метода можно найти в [54, 106].

Среди многошаговых методов наибольшее распространение в практике вычислений получили *методы Адамса*¹

$$\frac{y_{n+1} - y_n}{h} = \sum_{j=0}^k \beta_j f_{n+1-j}; \quad (14.75)$$

здесь $\beta_0, \beta_1, \dots, \beta_k$ — числовые коэффициенты, $f_{n+1-j} = f(t_{n+1-j}, y_{n+1-j})$.

Уравнение (14.75) позволяет найти новое значение y_{n+1} , используя найденные ранее значения $y_n, y_{n-1}, \dots, y_{n-k+1}$. Поэтому предварительно требуется задание k начальных значений y_0, y_1, \dots, y_{k-1} .

При $\beta_0 = 0$ метод Адамса является явным, так как значение y_{n+1} выражается через найденные ранее значения по явной формуле

$$y_{n+1} = y_n + h \sum_{j=1}^k \beta_j f_{n+1-j}. \quad (14.76)$$

Если же $\beta_0 \neq 0$, то для нахождения y_{n+1} приходится решать нелинейное уравнение

$$y_{n+1} = h\beta_0 f(t_{n+1}, y_{n+1}) + g_n, \quad (14.77)$$

где $g_n = y_n + h \sum_{j=1}^k \beta_j f_{n+1-j}$ — известное значение. Поэтому при $\beta_0 \neq 0$ метод Адамса (14.75) является неявным.

Существуют различные способы вывода формул (14.75). Приведем два из них. Воспользуемся, как и в §14.6, равенством

$$y(t_{n+1}) = y(t_n) + \int_{t_n}^{t_{n+1}} f(t, y(t)) dt. \quad (14.78)$$

Заменим приближенно функцию $F(t) \equiv f(t, y(t))$ интерполяционным многочленом $(k-1)$ -й степени $P_{k-1}(t)$, принимающим значения $f_{n+1-k}, \dots, f_{n-1}, f_n$ в тех узлах $t_{n+1-k}, \dots, t_{n-1}, t_n$, где значения сеточной функции y^h уже найдены. Интегрирование этого многочлена дает приближенное равенство

$$\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx \int_{t_n}^{t_{n+1}} P_{k-1}(t) dt = h \sum_{j=1}^k \beta_j f_{n+1-j}. \quad (14.79)$$

¹ Джон Кауч Адамс (1819—1892) — английский астроном и математик. Метод типа (14.75) был разработан им в 1855 г. по просьбе известного английского специалиста по внешней баллистике Башфорта.

В результате от (14.78) приходим к формуле (14.76), соответствующей явному *k*-шаговому методу Адамса—Башфорта.

Замечание. Так как многочлен P_{k-1} используется для приближения функции F вне отрезка, на котором известны ее значения, то в действительности равенство (14.79) основано на экстраполяции, поэтому соответствующий метод называют еще *экстраполяционным, методом Адамса*.

Если же в интеграле, входящем в равенство (14.78), заменить подынтегральную функцию интерполяционным многочленом *k*-й степени $Q_k(t)$, совпадающим со значениями $f_{n+1-k}, \dots, f_n, f_{n+1}$ в узлах $t_{n+1-k}, \dots, t_n, t_{n+1}$, то получится формула

$$y_{n+1} = y_n + h \sum_{j=0}^k \beta_j f_{n+1-j}, \quad (14.80)$$

соответствующая *k*-шаговому методу Адамса—Моултона. Заметим, что этот метод — неявный.

Замечание. Метод (14.80) принято называть также *интерполяционным методом Адамса*.

Выведем формулы двухшагового метода Адамса—Башфорта и одношагового метода Адамса—Моултона.

Интерполяционные многочлены $P_1(t)$ и $Q_1(t)$ таковы:

$$P_1(t) = f_{n-1}(t_n - t)/h + f_n(t - t_{n-1})/h,$$

$$Q_1(t) = f_n(t_{n+1} - t)/h + f_{n+1}(t - t_n)/h.$$

Их интегрирование по t дает следующие выражения:

$$\int_{t_n}^{t_{n+1}} P_1(t) dt = h \left(\frac{3}{2} f_n - \frac{1}{2} f_{n-1} \right), \quad \int_{t_n}^{t_{n+1}} Q_1(t) dt = h \left(\frac{1}{2} f_{n+1} + \frac{1}{2} f_n \right).$$

Таким образом, двухшаговая формула Адамса—Башфорта имеет вид

$$y_{n+1} = y_n + \frac{h}{2} (3f_n - f_{n-1}),$$

а одношаговая формула Адамса—Моултона — вид

$$y_{n+1} = y_n + \frac{h}{2} (f_{n+1} + f_n).$$

Предложение 14.1. Пусть решение задачи Коши $y(t)$ непрерывно дифференцируемо *k* раз на отрезке $[t_0, T]$. Тогда *k*-шаговый метод

Адамса—Башфорта и $(k - 1)$ -шаговый метод Адамса—Моултона имеют порядок аппроксимации, равный k . ■

Следующая теорема дает основание называть методы Адамса, имеющие p -й порядок аппроксимации, методами p -го порядка точности.

Теорема 14.8. *Пусть выполнено условие $|f'_y| \leq L$. Тогда явные методы Адамса устойчивы на конечном отрезке. Кроме того, при выполнении условия $h \leq h_0 = \frac{1}{2|\beta_0|L}$ устойчивыми на конечном отрезке являются и неявные методы Адамса.* ■

Замечание. Если выполнено условие $f'_y \leq 0$, то неявные методы Адамса устойчивы при любых h .

Следствие. *Пусть выполнено условие $|f'_y| \leq L$. Тогда если k -шаговый метод Адамса имеет p -й порядок аппроксимации, а начальные значения y_1, y_2, y_{k-1} определяются с p -м порядком точности, то метод сходится также с p -м порядком точности.*

Следствие верно в силу теоремы 14.4.

Приведем расчетные формулы методов Адамса—Башфорта p -го порядка точности при $p = 2, 3, 4$:

$$y_{n+1} = y_n + \frac{h}{2} (3f_n - f_{n-1}), \quad p = 2;$$

$$y_{n+1} = y_n + \frac{h}{12} (23f_n - 16f_{n-1} + 5f_{n-2}), \quad p = 3;$$

$$y_{n+1} = y_n + \frac{h}{24} (55f_n - 59f_{n-1} + 37f_{n-2} - 9f_{n-3}), \quad p = 4.$$

Приведем также расчетные формулы методов Адамса—Моултона p -го порядка точности при $p = 2, 3, 4$:

$$y_{n+1} = y_n + \frac{h}{2} (f_{n+1} + f_n), \quad p = 2,$$

$$y_{n+1} = y_n + \frac{h}{12} (5f_{n+1} + 8f_n - f_{n-1}), \quad p = 3;$$

$$y_{n+1} = y_n + \frac{h}{24} (9f_{n+1} + 19f_n - 5f_{n-1} + f_{n-2}), \quad p = 4.$$

2. Методы прогноза и коррекции. Может показаться, что при наличии явных формул Адамса высокого порядка точности нет необходимости в использовании неявных формул. Однако в вычислительной

практике явные методы Адамса используются очень редко. Одна из основных причин этого состоит в том, что в представляющих наибольший интерес для приложений задачах неявные методы обладают лучшими свойствами устойчивости и позволяют вести расчет с существенно большими шагами, нежели явные методы.

Сложность использования неявных методов Адамса заключается в необходимости решать уравнение (14.77) относительно y_{n+1} . Значение y_{n+1} можно найти, используя, например метод простой итерации

$$y_{n+1}^{(s+1)} = \psi(y_{n+1}^{(s)}), \quad s \geq 0; \quad \psi(y) = h\beta_0 f(t_{n+1}, y) + g_n. \quad (14.81)$$

Так как $\psi'(y) = h\beta_0 f'_y(t_{n+1}, y)$, то при достаточно малых h условие сходимости $|\psi'| \leq q < 1$ выполнено (см. § 4.4) и метод (14.81) сходится.

Часто за начальное приближение $y_{n+1}^{(0)}$ принимают значение, получаемое по явной формуле Адамса, и выполняют только одну итерацию метода (14.81). В результате приходят к *методу прогноза и коррекции*. Один из широко используемых методов прогноза и коррекции получается при совместном использовании методов Адамса—Башфорта и Адамса—Моултона четвертого порядка точности:

Прогноз:

$$y_{n+1}^{(0)} = y_n + \frac{h}{24} (55f_n - 59f_{n-1} + 37f_{n-2} - 9f_{n-3}),$$

$$f_{n+1}^{(0)} = f(t_{n+1}, y_{n+1}^{(0)}).$$

Коррекция:

$$y_{n+1} = y_n + \frac{h}{24} (9f_{n+1}^{(0)} + 19f_n - 5f_{n-1} + f_{n-2}).$$

Следует подчеркнуть, что результирующий метод оказался явным.

Пример 14.16. Применим описанный выше метод Адамса—Башфорта—Моултона четвертого порядка точности для решения задачи Коши (14.45) с шагом $h = 0.1$.

В качестве начальных значений y_1, y_2, y_3 необходимых для начала вычислений, примем значения, полученные методом Рунге—Кутты четвертого порядка точности и приведенные в табл. 14.4. Затем воспользуемся формулами

$$y_{n+1}^{(0)} = y_n + \frac{h}{12} (55t_n y_n - 59t_{n-1} y_{n-1} + 37t_{n-2} y_{n-2} - 9t_{n-3} y_{n-3}),$$

$$y_{n+1} = y_n + \frac{h}{12} (9t_{n+1} y_{n+1}^{(0)} + 19t_n y_n - 5t_{n-1} y_{n-1} + t_{n-2} y_{n-2}).$$

Таблица 14.5

t_n	Прогноз $y_n^{(0)}$	Погрешность прогноза $\varepsilon_n^{(0)}$	Коррекция y_n	Погрешность метода ε_n
0.4	1.173420048	$9.1 \cdot 10^{-5}$	1.173518429	$-7.6 \cdot 10^{-6}$
0.5	1.283880725	$1.5 \cdot 10^{-4}$	1.284044297	$-1.9 \cdot 10^{-5}$
0.6	1.433111448	$2.2 \cdot 10^{-4}$	1.433364614	$-3.5 \cdot 10^{-5}$
0.7	1.631994012	$3.3 \cdot 10^{-4}$	1.632374743	$-5.9 \cdot 10^{-5}$
0.8	1.896003688	$4.8 \cdot 10^{-4}$	1.896572568	$-9.2 \cdot 10^{-5}$
0.9	2.247194327	$7.2 \cdot 10^{-4}$	2.248046603	$-1.4 \cdot 10^{-4}$
1.0	2.717200091	$1.1 \cdot 10^{-3}$	2.718486351	$-2.1 \cdot 10^{-4}$

Найденные значения и соответствующие погрешности приведены в табл. 14.5. ▲

3. Общие линейные многошаговые методы. Эти методы, включающие в себя методы Адамса, задаются формулами вида

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y_{n+1-j} = \sum_{j=0}^k \beta_j f_{n+1-j}. \quad (14.82)$$

Предполагается, что $\alpha_0 \neq 0$, $|\alpha_k| + |\beta_k| \neq 0$. Они называются линейными, так как значения y_i и f_i ($i = n+1-k, \dots, n, n+1$) входят в формулу (14.82) линейно.

Замечание. Методы (14.82) принято также называть *конечно-разностными методами*, а дискретную задачу Коши для уравнения (14.82) — *конечно-разностной схемой* (или просто — *разностной схемой*).

4. Методы с переменным шагом и переменным порядком. На основе методов Адамса создан ряд весьма сложных, но и эффективных программ. В них предусматривается не только автоматический выбор шага (подобно тому, как это делается для методов Рунге—Кутты), но и автоматический выбор порядка метода. И шаг метода, и его порядок (в некоторых программах порядок точности может достигать 13) меняются в ходе вычислительного процесса, приспосабливаясь к характеру поведения искомого решения.

Методы Адамса требуют меньшего числа вычислений правой части дифференциального уравнения по сравнению с методами Рунге—Кутты того же порядка точности. Для них существуют эффективные

методы апостериорной оценки локальной погрешности. Недостатком методов Адамса является нестандартное начало вычислений. Для определения значений y_1, y_2, \dots, y_{k-1} , необходимых для работы k -шагового метода, используются методы Рунге—Кутты либо другие многошаговые методы. В разработанных к настоящему времени стандартных программах эта проблема решена.

§ 14.8. Устойчивость численных методов решения задачи Коши

Как в теории численных методов решения задачи Коши, так и в практическом плане вопросы об устойчивости методов к малым погрешностям задания начальных данных и правой части уравнения, а также об устойчивости к погрешностям вычислений являются одними из центральных. Рассмотрим некоторые естественные требования устойчивости, которые накладываются на дискретные методы

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y_{n+1-j} = \Phi(t_n, y_{n+1-k}, \dots, y_n, y_{n+1}, h). \quad (14.83)$$

Уделим сначала основное внимание исследованию устойчивости дискретной задачи Коши для уравнения (14.83) к малым погрешностям в начальных данных. Пусть y^h — решение дискретной задачи, соответствующее начальным значениям y_0, y_1, \dots, y_{k-1} , а y^{*h} — решение той же задачи, соответствующее начальным значениям $y_0^*, y_1^*, \dots, y_{k-1}^*$. Если отрезок $[t_0, T]$, на котором ищется решение, и шаг h фиксированы, то для всякого приемлемого метода решения задачи Коши его решение непрерывным образом зависит от начальных значений. Более того, если погрешности $\varepsilon_i = y_i - y_i^*$ ($i = 0, 1, \dots, k-1$) задания начальных данных достаточно малы, то погрешность значений y_n^* можно оценить следующим образом:

$$\max_{0 \leq n \leq N} |y_n^* - y_n| \leq K^h(T) \max_{0 \leq i \leq k-1} |\varepsilon_i|. \quad (14.84)$$

Величина $K^h(T)$, входящая в правую часть неравенства (14.84), играет роль числа обусловленности метода. Подчеркнем, что в общем случае она зависит как от T , так и от h .

1. Нуль-устойчивость. Будем стремиться к тому, чтобы при достаточно малых значениях шага h дискретная задача Коши не только имела близкое к $y(t)$ решение y^h , но и обладала другими важными свой-

ствами, аналогичными свойствам исходной задачи. В силу неравенства (14.11) погрешность, внесенная в начальное значение задачи Коши, на отрезке $[t_0, T]$ возрастает не более, чем в $K(T)$ раз (где $K(T)$, вообще говоря, растет с ростом T), поэтому в общем случае рост величин $K^h(T)$ с ростом T также допустим. Однако если коэффициент $K^h(T)$ может неограниченно возрастать при $h \rightarrow 0$, то уменьшение шага h приведет не к уточнению решения, а, наоборот, к неограниченному росту погрешности. Таким образом, следует потребовать, чтобы для дискретной задачи Коши при всех достаточно малых h было выполнено неравенство

$$\max_{0 \leq n \leq N} |y_n^* - y_n| \leq \bar{K}(T) \max_{0 \leq i \leq k-1} |\varepsilon_i|, \quad (14.85)$$

где $\bar{K}(T)$ не зависит от h .

Методы, для которых неравенство (14.85) выполнено в случае, когда решается задача Коши для однородного уравнения $y' = 0$, будем называть *нуль-устойчивыми*.

Чтобы отбросить те из методов (14.83), которые заведомо не обладают свойством нуль-устойчивости, применим метод (14.83) к решению задачи Коши для уравнения $y' = 0$. Тогда $f \equiv 0$, $\Phi \equiv 0$ и уравнение (14.83) принимает вид

$$\alpha_0 y_{n+1} + \alpha_1 y_n + \dots + \alpha_k y_{n-k+1} = 0. \quad (14.86)$$

Будем считать, что $\alpha_0 \neq 0$ и $\alpha_k \neq 0$. Такие уравнения называет *линейными однородными разностными уравнениями k -го порядка с постоянными коэффициентами*.

Пусть y^{*h} — решение того же уравнения, соответствующее возмущенным начальным значениям $y_0^*, y_1^*, \dots, y_{k-1}^*$. Тогда в силу линейности уравнения погрешность $\varepsilon_n = y_n - y_n^*$ также является его решением:

$$\alpha_0 \varepsilon_{n+1} + \alpha_1 \varepsilon_n + \dots + \alpha_k \varepsilon_{n-k+1} = 0. \quad (14.87)$$

Будем искать частное решение уравнения (14.87) в виде $\varepsilon_n = q^n$. Подставляя $\varepsilon_{n+1-j} = q^{n+1-j}$ ($j = 0, 1, \dots, k$) в (14.87), и сокращая на общий множитель q^{n+1-k} , видим, что величина q должна удовлетворять уравнению

$$\alpha_0 q^k + \alpha_1 q^{k-1} + \dots + \alpha_k = 0, \quad (14.88)$$

которое называют *характеристическим уравнением*, соответствующим методу (14.83). Многочлен $P(q) = \alpha_0 q^k + \alpha_1 q^{k-1} + \dots + \alpha_k$ называется *характеристическим многочленом*.

Приведем некоторые факты, известные из теории линейных разностных уравнений. Пусть q — корень уравнения (14.88) (вообще говоря, комплексный). Тогда сеточная функция $\varepsilon_n = q^n$ является решением разностного уравнения (14.87). Если же q — кратный корень кратности $m \geq 2$, то ему отвечают частные решения $q^n, nq^n, n^2q^n, \dots, n^{m-1}q^n$. Опишем теперь структуру общего решения разностного уравнения. Пусть q_1, q_2, \dots, q_r — корни характеристического уравнения, а m_1, m_2, \dots, m_r — их кратности ($m_1 + m_2 + \dots + m_r = k$). Тогда всякое решение уравнения (14.87) может быть представлено в виде¹

$$\varepsilon_n = \sum_{l=0}^{m_1-1} c_{1l} n^l q_1^n + \sum_{l=0}^{m_2-1} c_{2l} n^l q_2^n + \dots + \sum_{l=0}^{m_r-1} c_{rl} n^l q_r^n. \quad (14.89)$$

В частности, если все корни q_1, q_2, \dots, q_k — простые, то для всякого решения уравнения (14.87) справедливо представление

$$\varepsilon_n = c_1 q_1^n + c_2 q_2^n + \dots + c_k q_k^n.$$

Оказывается, что наличие или отсутствие у метода (14.83) нуль-устойчивости определяется исключительно расположением корней характеристического уравнения на комплексной плоскости.

Будем говорить, что выполнено *корневое условие*, если все корни q_1, q_2, \dots, q_r характеристического уравнения лежат внутри или на границе единичного круга комплексной плоскости (т.е. удовлетворяют условию $|q_i| \leq 1$, $i = 1, 2, \dots, r$), причем на границе единичного круга нет кратных корней. Заметим, что в силу равенства (14.30) число $q = 1$ всегда является корнем характеристического уравнения.

Теорема 14.9. Для того чтобы метод (14.83) обладал нуль-устойчивостью, необходимо и достаточно, чтобы выполнялось корневое условие.

Доказательство. Ограничимся доказательством необходимости выполнения корневого условия. Предположим, что метод обладает свойством нуль-устойчивости, а корневое условие не выполнено. Тогда характеристическое уравнение имеет либо корень q такой, что $|q| > 1$, либо корень q кратности $m \geq 2$ такой, что $|q| = 1$. В первом случае сеточная функция ε_n^h со значениями $\varepsilon_n = \varepsilon q^n$ есть решение разностного уравнения (14.87), соответствующее заданию начальных значений $\varepsilon_0 =$

¹ Если $\alpha_k = 0$ и среди корней характеристического уравнения имеется корень $q = 0$ кратности s , то представление (14.89) верно при $n \geq s$.

$= \varepsilon q^i$ ($i = 0, 1, \dots, k - 1$). Во втором случае решением, соответствующим заданию начальных значений $\varepsilon_i = \varepsilon i q^i$ ($i = 0, 1, \dots, k - 1$), является функция ε^h со значениями $\varepsilon_n = \varepsilon n q^n$. И в том, и в другом случаях благодаря выбору ε начальные погрешности могут быть сделаны сколь угодно малыми, но в то же время $|\varepsilon_n| \rightarrow \infty$ при $n \rightarrow \infty$.

Учитывая, что $N = (T - t_0)/h \rightarrow \infty$ при $h \rightarrow 0$, получаем $|y_N^* - y_N| = |\varepsilon_N| \rightarrow \infty$ при $h \rightarrow 0$, т.е. неравенство (14.85) не может выполняться для всех h . Итак, необходимость корневого условия доказана. ■

Теорема 14.10. *Методы Рунге—Кутты и Адамса обладают свойством нуль-устойчивости.*

Доказательство. Методам Рунге—Кутты соответствует однородное разностное уравнение $y_{n+1} - y_n = 0$, а ему, в свою очередь — характеристическое уравнение $q - 1 = 0$. Последнее имеет один простой корень $q = 1$, т.е. корневое условие выполнено.

Аналогично, k -шаговому методу Адамса соответствует разностное уравнение $y_{n+1} - y_n + 0 \cdot y_{n-1} + \dots + 0 \cdot y_{n-k+1} = 0$, а ему — характеристическое уравнение $q^k - q^{k-1} = 0$. Последнее имеет один простой корень $q_1 = 1$ и корень $q_2 = 0$ кратности $k - 1$, т.е. корневое условие здесь также выполняется. ■

Пример 14.17. Рассмотрим метод

$$\frac{y_{n+1} + y_n - 2y_{n-1}}{3h} = \frac{5f_n + f_{n-1}}{6}, \quad (14.90)$$

имеющий второй порядок аппроксимации¹. Попытаемся применить его для численного решения задачи Коши $y'(t) = \cos t$, $y(0) = 0$. Заметим, что функция $y(t) = \sin t$ является ее решением.

Возьмем шаг $h = 0.1$. Положим $y_0 = 0$ и в качестве второго начального значения, необходимого для расчета по методу (14.90), примем $y_1 = 0.1$. Так как $y(0.1) = 0.099833\dots$, то абсолютная погрешность значения y_1 не превышает $2 \cdot 10^{-4}$.

График полученного приближения изображен на рис. 14.13. При $t \geq 0.7$ погрешность уже становится заметной. Далее она быстро развивается и при $t \approx 1.4$ потеря точности становится катастрофической. Попытка увеличить точность решения за счет уменьшения шага вдвое приводит лишь к еще более быстрому нарастанию погрешности.

¹ В наличии второго порядка аппроксимации можно убедиться самостоятельно.

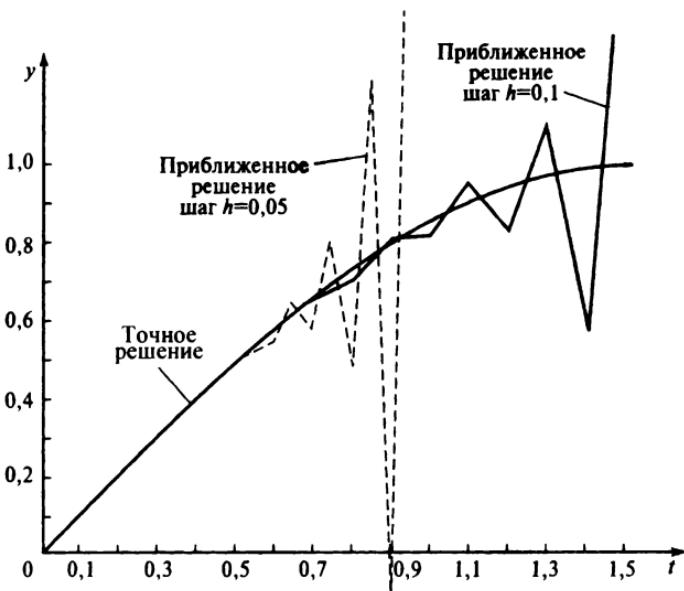


Рис. 14.13

Полная потеря точности происходит здесь уже при $t \approx 0.9$. Проверим теперь, выполняется ли для метода (14.90) корневое условие. Характеристическое уравнение имеет вид $q^2 + q - 2 = 0$. Корнями уравнения являются числа $q_1 = 1$, $q_2 = -2$. Так как $|q_2| = 2 > 1$, то корневое условие нарушено. Отсутствие у метода (14.90) нуль-устойчивости и служит причиной наблюдаемого неконтролируемого роста погрешности. Такого рода колебания приближенного решения, вызванные ростом погрешности, иногда называют *четно-нечетной болтанкой*¹. ▲

Таким образом, для того чтобы численный метод можно было использовать на практике, необходимо, чтобы он был нуль-устойчивым. Игнорирование этого требования даже для методов, обладающих высоким порядком аппроксимации, приводит к катастрофической потере точности.

Оказывается, что для линейных многошаговых методов выполнение корневого условия гарантирует не только нуль-устойчивость метода, но и устойчивость метода на конечном отрезке по начальным значениям и правой части в смысле определения устойчивости (14.28), введенного в § 14.2.

¹ Конечно, это не математический термин, а жаргон.

Теорема 14.11. Пусть выполнено условие $|f'_y| \leq L$. Предположим, что линейный многошаговый метод (14.82) удовлетворяет корневому условию и при $\beta_0 \neq 0$ (т.е. для неявного метода) выполнено дополнительное условие на шаг: $h \leq h_0 = \frac{|\alpha_0|}{2|\beta_0|L}$. Тогда метод (14.82) устойчив на конечном отрезке. ■

Доказательство этой теоремы можно найти в [88].

2. Абсолютная устойчивость. Как уже отмечалось в § 14.1, необходимость решения задачи Коши на больших временных отрезках $[t_0, T]$ возникает в самых различных областях науки и техники. Наибольший интерес при этом представляет изучение устойчивых решений. Отметим, что нуль-устойчивость гарантирует устойчивое развитие погрешностей при $h \rightarrow 0$ только в том случае, когда отрезок интегрирования $[t_0, T]$ фиксирован. Однако наличие нуль-устойчивости вовсе не исключает того, что сколь угодно малая погрешность в начальных значениях при неограниченном росте t_n (при $T \rightarrow \infty$) может приводить к сколь угодно большой погрешности решения. Значительную часть таких заведомо непригодных для решения задачи Коши на больших временных отрезках методов можно отбросить, если исследовать результат их применения к решению модельной задачи

$$y' = \lambda y, \quad y(t_0) = y_0. \quad (14.91)$$

Напомним (см. § 14.1), что решение этой задачи устойчиво по Ляпунову, если комплексный параметр λ удовлетворяет условию $\operatorname{Re} \lambda \leq 0$, и асимптотически устойчиво, если $\operatorname{Re} \lambda < 0$.

Большинство используемых дискретных методов (в том числе и методы Рунге—Кутты и Адамса) в применении к задаче (14.91) становятся линейными и приобретают вид

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y_{n+1-j} = \lambda \sum_{j=0}^k \beta_j (h\lambda) y_{n+1-j}; \quad (14.92)$$

здесь $\beta_j(h\lambda)$ — некоторые зависящие от величины $z = \lambda h$ функции.

В силу линейности уравнения (14.92) ошибка $\varepsilon_n = y_n - y_n^*$, возникающая из-за погрешностей в начальных значениях, удовлетворяет тому же уравнению:

$$\frac{1}{h} \sum_{j=0}^k \alpha_j \varepsilon_{n+1-j} = \lambda \sum_{j=0}^k \beta_j (h\lambda) \varepsilon_{n+1-j}.$$

Перепишем это уравнение в виде

$$\sum_{j=0}^k \gamma_j(z) \epsilon_{n+1-j} = 0, \quad (14.93)$$

где $\gamma_j(z) = \alpha_j - z\beta_j(z)$, $z = h\lambda$.

Заметим, что (14.93) — это линейное однородное разностное уравнение. Поэтому для того чтобы при фиксированном z погрешность ϵ_n оставалась ограниченной при $n \rightarrow \infty$, необходимо и достаточно выполнение корневого условия для отвечающего уравнению (14.93) полинома.

Назовем метод (14.83) *абсолютно устойчивым* для данного $z = h\lambda$, если при этом z все корни *полинома устойчивости*

$$P(q, z) = \sum_{j=0}^k \gamma_j(z) q^{k-j} \quad (14.94)$$

лежат в комплексной плоскости внутри единичного круга и на границе этого круга нет кратных корней. Множество D точек комплексной плоскости, состоящее из тех z , для которых метод абсолютно устойчив, называют *областью абсолютной устойчивости метода*.

Пример 14.18. Найдем область абсолютной устойчивости метода Эйлера.

Применимально к модельному уравнению $y' = \lambda y$ расчетная формула метода Эйлера принимает вид $y_{n+1} = y_n + h\lambda y_n$. Запишем соответствующее разностное уравнение для погрешности

$$\epsilon_{n+1} - (1 + h\lambda)\epsilon_n = 0 \quad (14.95)$$

и заметим, что полином устойчивости $P(q, z) = q - (1 + z)$ имеет один простой корень $q = 1 + z$. Область абсолютной устойчивости состоит здесь из тех $z = h\lambda$, для которых $|1 + z| \leq 1$, и представляет собой в комплексной плоскости круг единичного радиуса с центром в точке $z_0 = -1$ (рис. 14.14, а).

Тот же результат нетрудно установить, не используя полином устойчивости и корневое условие. Действительно, из уравнения (14.95) следует, что $\epsilon_n = (1 + h\lambda)^n \epsilon_0$. Поэтому $|\epsilon_n| = |1 + h\lambda|^n |\epsilon_0|$ и $\epsilon_n \rightarrow \infty$ при $n \rightarrow \infty$, если $|1 + h\lambda| > 1$. В этом случае метод не является устойчивым. Наоборот, если $|1 + h\lambda| \leq 1$, то $|\epsilon_n| \leq |\epsilon_0|$ и метод устойчив. ▲

Пример 14.19. Найдем области абсолютной устойчивости для неявного метода Эйлера (14.24), правила трапеций (14.25), метода Эйлера—Коши (14.61) и усовершенствованного метода Эйлера (14.62).

Применимально к модельному уравнению неявный метод Эйлера принимает вид $y_{n+1} = y_n + h\lambda y_{n+1}$. Соответствующий полином устой-

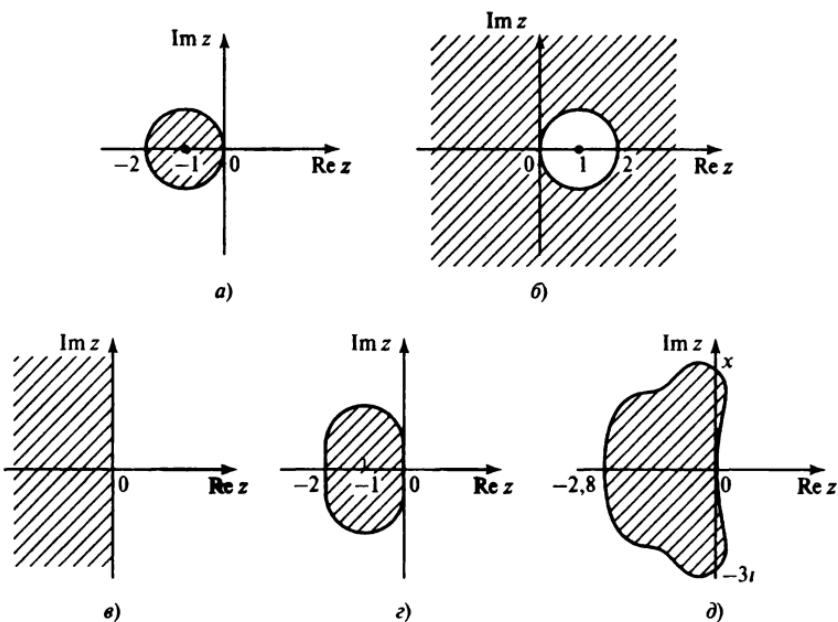


Рис. 14.14

чивости $(1 - z)q - 1$ имеет единственный корень $q = (1 - z)^{-1}$. Условие $|q| \leq 1$ эквивалентно здесь условию $|1 - z| \geq 1$. Таким образом, область устойчивости представляет собой внешнюю часть единичного круга с центром в точке $z_0 = 1$ (рис. 14.14, б).

Для правила трапеций $y_{n+1} = y_n + \frac{h\lambda}{2}(y_n + y_{n+1})$ полином устойчивости $P(q, z) = \left(1 - \frac{z}{2}\right)q - \left(1 + \frac{z}{2}\right)$ имеет корень $q = \frac{2+z}{2-z}$. Условие $|q| \leq 1$ эквивалентно здесь неравенству $\text{Re } z \leq 0$, поэтому область устойчивости представляет собой левую полуплоскость (рис. 14.14, в).

Применимтельно к уравнению $y' = \lambda y$ расчетные формулы методов (14.61) и (14.62) совпадают: $y_{n+1} = y_n + \frac{h\lambda}{2}(2 + h\lambda)y_n$. Корнем полинома устойчивости $q - \left(1 + z + \frac{z^2}{2}\right)$ является $q = 1 + z + \frac{z^2}{2}$. Область абсолютной устойчивости изображена на рис. 14.14, г. Для сравнения на рис. 14.14, д схематично изображена область абсолютной устойчивости метода Рунге—Кутты четвертого порядка точности. ▲

Предположим, что параметр λ , входящий в модельное уравнение (14.91), отрицателен. Тогда условие абсолютной устойчивости метода Эйлера $|1 + h\lambda| \leq 1$ оказывается эквивалентным неравенству

$$h \leq h_0 = \frac{2}{|\lambda|}. \quad (14.96)$$

Такое же ограничение на шаг возникает при использовании метода Эйлера—Коши и усовершенствованного метода Эйлера. В то же время метод Рунге—Кутты четвертого порядка точности оказывается абсолютно устойчивым при выполнении чуть менее ограничительного условия

$$h \leq h_0 \approx \frac{2.8}{|\lambda|}. \quad (14.97)$$

Отметим, что при $\lambda < 0$ неявный метод Эйлера и правило трапеций оказываются абсолютно устойчивыми при любых h .

3. A -устойчивость. Для того чтобы исключить ограничение на шаг h при решении устойчивой по Ляпунову модельной задачи (14.91), необходимо потребовать, чтобы область абсолютной устойчивости метода включала в себя полуплоскость $\operatorname{Re} z < 0$. Численный метод, обладающий таким свойством, называют *A-устойчивым*.

Примерами *A*-устойчивых методов служат неявный метод Эйлера и правило трапеций. В то же время метод Эйлера и метод Рунге—Кутты четвертого порядка точности не являются *A*-устойчивыми.

§ 14.9. Неявный метод Эйлера

Как следует из результатов предыдущего параграфа, простейшим представителем семейства *A*-устойчивых методов является неявный метод Эйлера

$$y_{n+1} = y_n + hf(t_{n+1}, y_{n+1}). \quad (14.98)$$

Как нетрудно понять, геометрическая интерпретация одного шага метода (14.98) заключается в том, что решение на отрезке $[t_n, t_{n+1}]$ аппроксимируется касательной $y = y_{n+1} + y'(t_{n+1})(t - t_{n+1})$, проведенной в точке (t_{n+1}, y_{n+1}) к интегральной кривой, проходящей через эту точку (рис. 14.15).

2. Устойчивость. Достоинства неявного метода Эйлера проявляются при решении дифференциальных уравнений, имеющих устойчивые по Ляпунову решения. Как уже отмечалось в § 14.2, достаточным условием такой устойчивости является выполнение одностороннего условия Липшица $f'_y \leq 0$.

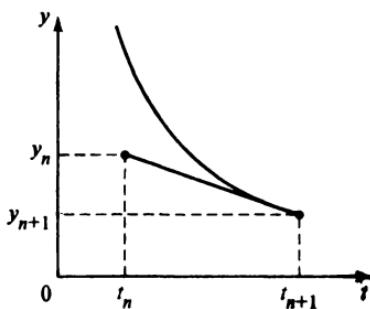


Рис. 14.15

Пусть y^h — решение дискретной задачи Коши для уравнения (14.98), соответствующее начальному условию $y^h(t_0) = y_0$, а y^{*h} — решение возмущенной задачи

$$y_{n+1}^* = y_n^* + h(f(t_{n+1}, y_{n+1}^*) + \psi_n), \quad (14.99)$$

$$y^{*h}(t_0) = y_0^*.$$

Аналогично теореме 14.5 устанавливается следующий результат.

Теорема 14.12. Пусть функция f удовлетворяет условию $f'_y \leq 0$. Тогда справедливо неравенство

$$\max_{0 \leq n \leq N} |y_n^* - y_n| \leq |y_0^* - y_0| + h \sum_{k=0}^{N-1} |\psi_k|, \quad (14.100)$$

означающее, что неявный метод Эйлера устойчив на конечном отрезке.

Доказательство. Вычитая из равенства (14.98) равенство (14.99), получаем

$$y_{n+1}^* - y_{n+1} = y_n^* - y_n + h f'_y(t_{n+1}, \tilde{y}_{n+1})(y_{n+1}^* - y_{n+1}) + h \psi_n,$$

где \tilde{y}_{n+1} — некоторое промежуточное между y_{n+1} и y_{n+1}^* значение. Откуда

$$[1 - h f'_y(t_{n+1}, \tilde{y}_{n+1})] |y_{n+1}^* - y_{n+1}| \leq |y_n^* - y_n| + h |\psi_n|. \quad (14.101)$$

Поскольку $f'_y \leq 0$, то справедливо неравенство

$$|y_{n+1}^* - y_{n+1}| \leq |y_n^* - y_n| + h |\psi_n|, \quad 0 \leq n < N,$$

из которого вытекает оценка

$$|y_n^* - y_n| \leq |y_0^* - y_0| + h \sum_{k=0}^{n-1} |\psi_k|, \quad 1 \leq n \leq N,$$

и, как следствие, — неравенство (14.100). ■

Замечание. Неравенство (14.100) является дискретным аналогом оценки (14.16), справедливой для погрешности решения задачи Коши (при этом $K(T) = 1$). Следует отметить, что оно верно для метода (14.98) при любых h , в то время как для явного метода Эйлера это неравенство имеет место только если $-L \leq f'_v \leq 0$ и $h \leq h_0 = 2/L$.

Для неявного метода Эйлера справедлив и дискретный аналог оценки (14.17).

Теорема 14.13. Пусть функция f удовлетворяет условию $f'_y \leq \sigma < 0$. Тогда справедливо неравенство

$$|y_n^* - y_n| \leq e^{\sigma_h(\tau_n - \tau_0)} |y_0^* - y_0| + \frac{1}{|\sigma|} \max_{0 \leq k < n} |\psi_k|, \quad n \geq 1, \quad (14.102)$$

где $\sigma_h = \frac{\sigma}{1 - \sigma h} \rightarrow \sigma$ при $h \rightarrow 0$.

Доказательство. Учитывая, что $f'_y \leq \sigma < 0$, из неравенства (14.101) выводим

$$|y_{n+1}^* - y_{n+1}| \leq \frac{1}{1 - h\sigma} |y_n^* - y_n| + \frac{h}{1 - h\sigma} |\psi_n|. \quad (14.103)$$

Докажем теперь оценку (14.102) используя метод математической индукции.

При $n = 1$ эта оценка имеет вид

$$|y_1^* - y_1| \leq e^{\sigma_h h} |y_0^* - y_0| + \frac{1}{|\sigma|} |\psi_0|$$

и следует из (14.103) с учетом неравенства $\frac{1}{1 + \alpha} \leq e^{-\alpha/(1 + \alpha)}$ с $\alpha = -h\sigma$

и оценки $\frac{h}{1 - h\sigma} \leq \frac{1}{|\sigma|}$.

Пусть оценка (14.102) верна при $n = k$. Тогда

$$\begin{aligned} |y_{k+1}^* - y_{k+1}| &\leq \frac{1}{1-h\sigma} |y_k^* - y_k| + \frac{h}{1-h\sigma} |\psi_k| \leq e^{\sigma h} e^{\sigma h(t_k - t_0)} |y_0^* - y_0| + \\ &+ \frac{1}{1-h\sigma} \left(\frac{1}{|\sigma|} + h \right) \max_{0 \leq i \leq k} |\psi_i| = e^{\sigma h(t_{k+1} - t_0)} |y_0^* - y_0| + \frac{1}{|\sigma|} \max_{0 \leq i \leq k} |\psi_i|. \end{aligned}$$

Следовательно оценка (14.102) верна и при $n = k + 1$. ■

3. Оценка погрешности. Из того, что неявный метод Эйлера устойчив и имеет первый порядок аппроксимации, вытекает (в силу теоремы 14.4) его сходимость с первым порядком. Приведем соответствующий результат.

Теорема 14.14. Пусть функция f удовлетворяет условию $f'_y \leq 0$. Тогда для неявного метода Эйлера справедлива следующая оценка погрешности:

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq \int_0^T |y''(t)| dt \cdot h. \quad (14.104)$$

Если же функция f удовлетворяет условию $f'_y \leq \sigma < 0$, то верна оценка

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq \frac{M_2}{2|\sigma|} h, \quad (14.105)$$

где $M_2 = \max_{[t_0, T]} |y''(t)|$.

Доказательство. Для погрешности аппроксимации ψ_n имеем

$$\begin{aligned} \Psi_n &= \frac{y(t_{n+1}) - y(t_n)}{h} - f(t_{n+1}, y(t_{n+1})) = \frac{y(t_{n+1}) - y(t_n)}{h} - y'(t_{n+1}) = \\ &= \frac{1}{h} \int_{t_n}^{t_{n+1}} (y'(s) - y'(t_{n+1})) ds = - \frac{1}{h} \int_{t_n}^{t_{n+1}} \left(\int_s^{t_{n+1}} y''(t) dt \right) ds = - \int_{t_n}^{t_{n+1}} \frac{t - t_n}{h} y''(t) dt. \end{aligned}$$

Таким образом, значения $y_n^* = y(t_n)$ удовлетворяют уравнениям (14.99),

в которых $|\psi_n| \leq \int_{t_n}^{t_{n+1}} |y''(t)| dt$. Применяя теорему 14.12, приходим к оценке (14.104).

Если же использовать оценку $|\psi_n| \leq \frac{1}{2} M_2 h$ и применить теорему 14.13, то результатом будет неравенство (14.105). ■

Замечание. Оценка (14.105) замечательна тем, что ее правая часть не растет с ростом T , если вторая производная решения ограничена.

4. Влияние погрешности вычислений. Поскольку метод (14.98) — неявный, то при его реализации на компьютере будет использован некоторый итерационный метод и будут найдены приближенные значения решения y_n^* удовлетворяющие соотношениям

$$y_{n+1}^* = y_n^* + hf(t_{n+1}, y_{n+1}^*) + \delta_n.$$

Величины δ_n — это невязки, с которыми найденные значения удовлетворяют уравнениям (14.98).

Предположим, что функция f удовлетворяет условию $f'_y \leq \sigma < 0$. Тогда, полагая $\psi_n = \delta_n/h$ и используя теорему 14.13, имеем

$$\max_{0 \leq n \leq N} |y_n^* - y_n| \leq \frac{1}{|\sigma| h} \max_{0 \leq n < N} |\delta_n|.$$

Заметим, что в этой оценке (в отличие от соответствующей оценки (14.54) для метода Эйлера) отсутствует множитель, растущий с ростом T .

§ 14.10. Решение задачи Коши для систем обыкновенных дифференциальных уравнений и дифференциальных уравнений m -го порядка

1. Задача Коши для систем дифференциальных уравнений первого порядка. Как правило, возникающие на практике проблемы приводят к необходимости решать задачу Коши не для одного дифференциального уравнения, а для систем дифференциальных уравнений вида

$$\begin{aligned} y'_1(t) &= f_1(t, y_1(t), y_2(t), \dots, y_m(t)), \\ y'_2(t) &= f_2(t, y_1(t), y_2(t), \dots, y_m(t)), \\ &\dots \\ y'_m(t) &= f_m(t, y_1(t), y_2(t), \dots, y_m(t)), \end{aligned} \tag{14.106}$$

где $y_1(t), y_2(t), \dots, y_m(t)$ — искомые функции, значения которых подлежат определению при $t \in [t_0, T]$.

В момент времени $t = t_0$ задаются *начальные условия*

$$y_1(t_0) = y_{10}, \quad y_2(t_0) = y_{20}, \dots, \quad y_m(t_0) = y_{m0}, \quad (14.107)$$

определяющие начальное состояние физической системы, развитие которой описывается уравнениями (14.106).

Введем вектор-функции $y(t) = (y_1(t), y_2(t), \dots, y_m(t))^T$, $f(t, y) = (f_1(t, y), f_2(t, y), \dots, f_m(t, y))^T$ и вектор $y_0 = (y_{10}, y_{20}, \dots, y_{m0})^T$. Тогда задачу Коши (14.106), (14.107) можно записать в компактной форме:

$$y'(t) = f(t, y(t)), \quad (14.108)$$

$$y(t_0) = y_0. \quad (14.109)$$

Для того чтобы охватить ряд важных для технических областей задач (электротехника, радиотехника и др.), будем считать, что функции y_i и f_i ($i = 1, 2, \dots, m$) могут принимать комплексные значения.

2. Разрешимость задачи Коши. Пусть Π_T — множество таких точек (t, y) , для которых $t \in [t_0, T]$, а y_1, y_2, \dots, y_m — произвольные комплексные числа. Это множество будем называть *слоем*. Будем, как и ранее, использовать обозначения $(y, z) = \sum_{i=1}^m y_i \bar{z}_i$ и $\|y\| = \left(\sum_{i=1}^m |y_i|^2 \right)^{1/2}$

для скалярного произведения и евклидовой нормы m -мерных комплексных векторов.

Сформулируем аналог теоремы 14.1 о разрешимости задачи Коши.

Теорема 14.15. Пусть вектор-функция $f(t, y)$ определена и непрерывна в слое Π_T . Предположил также, что она удовлетворяет условию Липшица

$$\|f(t, y_1) - f(t, y_2)\| \leq L \|y_1 - y_2\| \quad (14.110)$$

для всех $t_0 \leq t \leq T$ и произвольных y_1, y_2 , где $L > 0$ — некоторая постоянная (постоянная Липшица).

Тогда для каждого начального значения y_0 существует единственное решение $y(t)$ задачи Коши (14.108), (14.109), определенное на отрезке $[t_0, T]$. ■

Замечание 1. Можно показать, что если функции f_1, f_2, \dots, f_m непрерывно дифференцируемы по y_1, y_2, \dots, y_m , то условие Липшица (14.110) выполняется с постоянной L тогда и только тогда, когда матрица Якоби

$$f'_y(t, y) = \begin{bmatrix} f'_{1y_1}(t, y) & f'_{1y_2}(t, y) & \dots & f'_{1y_m}(t, y) \\ f'_{2y_1}(t, y) & f'_{2y_2}(t, y) & \dots & f'_{2y_m}(t, y) \\ \dots & \dots & \dots & \dots \\ f'_{my_1}(t, y) & f'_{my_2}(t, y) & \dots & f'_{my_m}(t, y) \end{bmatrix}$$

удовлетворяет неравенству $\|f'_y(t, y)\| \leq L$.

Замечание 2. Теорема 14.15 остается справедливой, если в ее формулировке условие Липшица (14.110) заменить менее ограничительным *односторонним условием Липшица*

$$\operatorname{Re}(f(t, y_1) - f(t, y_2), y_1 - y_2) \leq \sigma \|y_1 - y_2\|^2. \quad (14.111)$$

Назовем систему дифференциальных уравнений *диссипативной* если вектор-функция удовлетворяет неравенству

$$\operatorname{Re}(f(t, y_1) - f(t, y_2), y_1 - y_2) \leq 0$$

(т.е. если f удовлетворяет одностороннему условию Липшица с постоянной $\sigma = 0$).

3. Устойчивость решения задачи Коши. Приведем аналог теоремы 14.3.

Теорема 14.16. Пусть выполнены условия теоремы 14.15. Далее, пусть $y(t)$ — решение задачи (14.108), (14.109), а $y^*(t)$ — решение задачи

$$(y^*)'(t) = f(t, y^*(t)) + \psi(t), \quad (14.112)$$

$$y^*(t_0) = y_0^*. \quad (14.113)$$

Тогда справедлива оценка

$$\max_{t_0 \leq t \leq T} \|y(t) - y^*(t)\| \leq K(T) \left(\|y_0 - y_0^*\| + \int_{t_0}^T \|\psi(t)\| dt \right), \quad (14.114)$$

выражающая устойчивость на конечном отрезке $[t_0, T]$ решения задачи Коши по начальным значениям и правой части. Здесь $K(T) = e^{L(T-t_0)}$. ■

Если в теореме 14.15 условие Липшица (14.110) заменить односторонним условием (14.111), то оценка (14.114) будет выполнена с постоянной $K(T) = e^{\sigma(T-t_0)}$ при $\sigma > 0$ и с постоянной $K(T) = 1$ при $\sigma \leq 0$.

Следствие. Если система (14.108) диссипативна, то справедлива оценка

$$\max_{t_0 \leq t \leq T} \|y(t) - y^*(t)\| \leq \|y_0 - y_0^*\| + \int_{t_0}^T \|x(t)\| dt. \quad (14.115)$$

Замечание. Можно показать, что если условие (14.111) выполнено с постоянной $\sigma < 0$, то для всех $t \in [t_0, T]$ справедлива оценка

$$\|y(t) - y^*(t)\| \leq e^{\sigma(t-t_0)} \|y_0 - y_0^*\| + \frac{1}{|\sigma|} \max_{t_0 \leq t' \leq t} \|x(t')\|. \quad (14.116)$$

По аналогии со случаем одного дифференциального уравнения (см. § 14.1) рассмотрим вопрос об устойчивости решения задачи Коши к возмущениям начальных данных при $T \rightarrow \infty$. Будем считать, что на каждом отрезке $[t_0, T]$ ($t_0 < T$ — произвольно) неравенство (14.111) выполнено с некоторой постоянной $\sigma = \sigma(T)$. Тогда решение $y(t)$ определено для всех $t_0 \leq t < \infty$. Пусть $y^*(t)$ — решение задачи (14.112), (14.113), отвечающее произвольному начальному значению y_0^* и $\psi(t) \equiv 0$. Будем называть решение задачи Коши (14.108), (14.109) *устойчивым по Ляпунову*, если справедлива оценка $\max_{t_0 \leq t' \leq T} \|y(t) - y^*(t)\| \leq K \|y_0 - y_0^*\|$, где постоянная K не зависит от T . Если дополнительно известно, что $\|y(t) - y^*(t)\| \rightarrow 0$ при $t \rightarrow \infty$, то решение называется *асимптотически устойчивым*.

Замечание 1. При $\psi(t) \equiv 0$ из неравенства (14.115) следует оценка $\max_{t_0 \leq t \leq T} \|y(t) - y^*(t)\| \leq \|y_0 - y_0^*\|$, поэтому всякое решение диссипативной системы является устойчивым по Ляпунову.

Замечание 2. Предположим, что одностороннее условие Липшица (14.111) для всех $t_0 \leq t < \infty$ выполняется с одной и той же постоянной $\sigma < 0$. Тогда при $\psi(t) \equiv 0$ из неравенства (14.116) получаем оценку

$$\|y(t) - y^*(t)\| \leq e^{\sigma(t-t_0)} \|y_0 - y_0^*\|,$$

откуда следует, что $\|y(t) - y^*(t)\| \rightarrow 0$, при $t \rightarrow \infty$, т.е. решение $y(t)$ асимптотически устойчиво.

4. Система линейных дифференциальных уравнений с постоянными коэффициентами. Рассмотрим систему

$$\mathbf{y}'(t) = \mathbf{A}\mathbf{y}(t), \quad (14.117)$$

являющуюся простейшим примером системы дифференциальных уравнений первого порядка; здесь \mathbf{A} — квадратная матрица порядка m . В теории численных методов решения систем дифференциальных уравнений система (14.117) играет роль, аналогичную той, которую при исследовании методов численного интегрирования одного уравнения $y' = f(t, y)$ выполняет модельное уравнение $y'(t) = \lambda y(t)$ (см. § 14.1).

Напомним структуру решения системы линейных уравнений с постоянными коэффициентами в наиболее простом и важном случае, когда матрица \mathbf{A} имеет простую структуру. В этом случае существует набор $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_m$ собственных векторов матрицы \mathbf{A} , соответствующих собственным значениям $\lambda_1, \lambda_2, \dots, \lambda_m$, который образует базис в пространстве m -мерных векторов. Матрица \mathbf{P} , столбцами которой служат векторы $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_m$, не вырождена и такова, что

$$\mathbf{P}^{-1}\mathbf{A}\mathbf{P} = \Lambda = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \lambda_m \end{bmatrix}. \quad (14.118)$$

Обозначим через $z_1(t), z_2(t), \dots, z_m(t)$ координаты вектора $\mathbf{y}(t)$ в базисе $\mathbf{e}_1, \dots, \mathbf{e}_m$. Так как

$$\mathbf{y}(t) = z_1(t)\mathbf{e}_1 + z_2(t)\mathbf{e}_2 + \dots + z_m(t)\mathbf{e}_m, \quad (14.119)$$

то вектор $\mathbf{y}(t)$ связан с вектором $\mathbf{z}(t) = (z_1(t), z_2(t), \dots, z_m(t))^T$ равенством $\mathbf{y}(t) = \mathbf{P}\mathbf{z}(t)$. Умножив обе части системы (14.117) слева на \mathbf{P}^{-1} , получим для вектор-функции $\mathbf{z}(t)$ систему уравнений

$$\mathbf{z}'(t) = \Lambda\mathbf{z}(t).$$

В силу диагональной структуры матрицы Λ эта система распадается на m независимых дифференциальных уравнений

$$z'_i(t) = \lambda_i z_i(t), \quad i = 1, 2, \dots, m. \quad (14.120)$$

Заметим, что уравнение (14.120) для i -й компоненты вектора \mathbf{z} есть модельное уравнение с параметром $\lambda = \lambda_i$.

Таким образом, в рассматриваемом случае интегрирование системы линейных дифференциальных уравнений эквивалентно интегрирова-

нию m модельных уравнений (14.120). Решая каждое из них, получаем $z_i(t) = c_i e^{\lambda_i(t-t_0)}$, $c_i = z_i(t_0)$ и в силу равенства (14.119) находим

$$y(t) = \sum_{i=1}^m c_i e^{\lambda_i(t-t_0)} e_i.$$

Здесь $\mathbf{c} = (c_1, c_2, \dots, c_m)^T = \mathbf{P}^{-1} \mathbf{y}_0$.

Пусть $\varepsilon(t) = y(t) - y^*(t)$ — погрешность решения, вызванная погрешностью начальных значений $\varepsilon_0 = \mathbf{y}_0 - \mathbf{y}_0^*$. В силу линейности системы (14.119) погрешность является решением той же системы: $\varepsilon'(t) = A\varepsilon(t)$. Следовательно,

$$\mathbf{e}(t) = \sum_{i=1}^m \alpha_i e^{\lambda_i(t-t_0)} e_i, \quad (14.121)$$

где $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_m)^T = \mathbf{P}^{-1} \varepsilon_0$.

Формула (14.121) позволяет сделать ряд важных выводов. В частности, из нее следует, что решение системы (14.117) с постоянной матрицей A простой структуры устойчиво по Ляпунову тогда и только тогда, когда $\operatorname{Re} \lambda_i \leq 0$ для всех $i = 1, 2, \dots, m$ и асимптотически устойчиво тогда и только тогда, когда $\operatorname{Re} \lambda_i < 0$ для всех $i = 1, 2, \dots, m$. Отметим также, что в отличие от случая одного дифференциального уравнения для систем характерно наличие m временных постоянных $\tau_i = \frac{1}{|\operatorname{Re} \lambda_i|}$,

$i = 1, 2, \dots, m$. Наличие в решении задачи физических компонент с существенно различными временными постоянными может привести к серьезным затруднениям при численном решении соответствующих задач. Более подробно этот вопрос рассматривается в § 14.11.

Пусть теперь $y^*(t)$ — решение нелинейной системы

$$(y^*)'(t) = f(t, y^*(t)), \quad (14.122)$$

отвечающие возмущенному начальному условию $y^*(t_0) = \mathbf{y}_0^*$. Вычитая из уравнения $y'(t) = f(t, y(t))$ уравнение (14.122) и используя приближенное равенство

$$f(t, y(t)) - f(t, y^*(t)) \approx f'_y(t, y(t))(y(t) - y^*(t)),$$

получаем, что погрешность $\varepsilon(t) = y(t) - y^*(t)$ удовлетворяет приближенному равенству

$$\varepsilon'(t) \approx A\varepsilon(t),$$

где $A = f'_y(t, y(t))$.

Таким образом, можно предположить, что в малой окрестности точки $(\tilde{t}, y(\tilde{t}))$ эволюция погрешности, $\varepsilon(t)$ происходит примерно так, как и для системы (14.119), т.е.

$$\varepsilon(t) \approx \sum_{i=1}^m \alpha_i e^{\lambda_i(t-\tilde{t})} e_i.$$

Здесь λ_i и e_i ($i = 1, \dots, m$) — собственные значения и собственные векторы матрицы $A = f'_y(\tilde{t}, y(\tilde{t}))$.

5. Понятие о численных методах решения задачи Коши для систем уравнений первого порядка. Описанные выше применительно к решению задачи Коши для одного уравнения методы можно использовать и для систем уравнений первого порядка, причем форма их записи претерпевает минимальные изменения. Следует лишь заменить в расчетных формулах числа y_n на векторы $y_n = (y_{1n}, y_{2n}, \dots, y_{mn})^T$, функцию f — на вектор-функцию f и т.д. В результате дискретное уравнение (14.18) преобразуется в систему дискретных уравнений

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y_{n+1-j} = \Phi(t_n, y_{n+1-k}, \dots, y_n, y_{n+1}, h).$$

Например, расчетная формула метода Эйлера $y_{n+1} = y_n + hf(t_n, y_n)$ применительно к решению системы (14.108) принимает вид

$$y_{n+1} = y_n + hf(t_n, y_n).$$

Покоординатная запись этого соотношения выглядит так:

$$y_{1,n+1} = y_{1n} + hf_1(t_n, y_{1n}, y_{2n}, \dots, y_{mn}),$$

$$y_{2,n+1} = y_{2n} + hf_2(t_n, y_{1n}, y_{2n}, \dots, y_{mn}),$$

.....

$$y_{m,n+1} = y_{mn} + hf_m(t_n, y_{1n}, y_{2n}, \dots, y_{mn}).$$

Аналогично, метод Рунге—Кутты четвертого порядка точности (14.70) порождает для систем дифференциальных уравнений первого порядка следующий метод:

$$y_{n+1} = y_n + hk_n, \quad k_n = \frac{1}{6} (k_n^{(1)} + 2k_n^{(2)} + 2k_n^{(3)} + k_n^{(4)});$$

$$k_n^{(1)} = f(t_n, y_n), \quad k_n^{(2)} = f\left(t_n + \frac{h}{2}, y_n + \frac{h}{2} k_n^{(1)}\right),$$

$$k_n^{(3)} = f\left(t_n + \frac{h}{2}, y_n + \frac{h}{2} k_n^{(2)}\right), \quad k_n^{(4)} = f(t_n + h, y_n + hk_n^{(3)}).$$

Теория численных методов решения задачи Коши для систем дифференциальных уравнений имеет много общего с соответствующей теорией решения задачи Коши для одного дифференциального уравнения. В частности, справедливы аналоги всех изложенных выше результатов, касающихся устойчивости и сходимости дискретных методов на конечном отрезке. Однако имеют место и существенно новые явления. Один из таких эффектов — жесткость — будет рассмотрен в § 14.11. Прежде чем переходить к его изложению, выясним, какие изменения появляются при применении дискретных методов к решению задачи Коши для системы уравнений с постоянными коэффициентами (14.117).

Рассмотрим линейный многошаговый метод

$$\frac{1}{h} \sum_{j=0}^k \alpha_j y_{n+1-j} = \sum_{j=0}^k \beta_j A y_{n+1-j}. \quad (14.123)$$

Предположим для упрощения, что матрица A имеет простую структуру.

Положим $\mathbf{z}_n = P^{-1} y_n$, где P — матрица, удовлетворяющая равенству (14.118). Умножая обе части уравнения (14.123) на матрицу P^{-1} слева и учитывая, что

$$y_{n+1-j} = P \mathbf{z}_{n+1-j}, \quad P^{-1} A y_{n+1-j} = P^{-1} A P \mathbf{z}_{n+1-j} = \Lambda \mathbf{z}_{n+1-j},$$

получаем соотношение

$$\frac{1}{h} \sum_{j=0}^k \alpha_j \mathbf{z}_{n+1-j} = \sum_{j=0}^k \beta_j \Lambda \mathbf{z}_{n+1-j}.$$

Так как матрица Λ — диагональная, то оно эквивалентно следующей системе уравнений для компонент вектора $\mathbf{z}_n = (z_{1n}, z_{2n}, \dots, z_{mn})^T$:

$$\frac{1}{h} \sum_{j=0}^k \alpha_j z_{i,n+1-j} = \sum_{j=0}^k \beta_j \lambda_i z_{i,n+1-j}, \quad i = 1, 2, \dots, m. \quad (14.124)$$

Заметим, что (14.124) есть не что иное, как результат применения линейного многошагового метода к решению уравнений (14.120).

Таким образом, если решения системы (14.117) устойчивы по Ляпунову, то для того чтобы погрешности ε_n оставались ограниченными при $n \rightarrow \infty$, необходимо потребовать, чтобы для всех $i = 1, \dots, m$ величина $h\lambda_i$ принадлежала области D абсолютной устойчивости применяемого метода. Можно показать, что такой же вывод справедлив и для методов Рунге—Кутты. Требование, чтобы $h\lambda_i \in D$ при всех $i = 1, 2, \dots, m$, для методов, не обладающих свойством A -устойчивости, может приводить к существенным ограничениям на длину шага h .

Предположим, например, что все собственные значения матрицы A отрицательны. Тогда условие (14.96) абсолютной устойчивости

метода Эйлера приводит к следующему ограничению на длину шага интегрирования:

$$h \leq h_0 = \min_{1 \leq i \leq m} \frac{2}{|\lambda_i|}. \quad (14.125)$$

Такое же ограничение на шаг возникает при использовании метода Эйлера—Коши и усовершенствованного метода Эйлера. Метод Рунге—Кутты четвертого порядка точности, как следует из неравенства (14.97), оказывается абсолютно устойчивым при таком ограничении на длину шага:

$$h \leq h_0 \approx \min_{1 \leq i \leq m} \frac{2.8}{|\lambda_i|}.$$

Следовательно, для явных методов шаг интегрирования должен не превышать значения h_0 , пропорционального наименьшей из временных постоянных системы.

6. Сведение задачи Коши для уравнения m -го порядка к задаче Коши для системы уравнений первого порядка. Задача Коши для дифференциального уравнения m -го порядка состоит в нахождении функции $y(t)$, удовлетворяющей при $t \geq t_0$ дифференциальному уравнению

$$y^{(m)}(t) = f(t, y(t), y'(t), \dots, y^{(m-1)}(t)), \quad (14.126)$$

а при $t = t_0$ — начальным условиям

$$y(t_0) = y_{10}, \quad y'(t_0) = y_{20}, \dots, \quad y^{(m-1)}(t_0) = y_{m0}. \quad (14.127)$$

Рассмотрим функции $y_1(t) = y(t)$, $y_2(t) = y'(t)$, ..., $y_m(t) = y^{(m-1)}(t)$. Заметим, что $y_k(t) = y'_{k-1}(t)$. Поэтому введенные функции удовлетворяют системе дифференциальных уравнений первого порядка

$$\begin{aligned} y'_1 &= y_2, \\ y'_2 &= y_3, \\ &\dots \dots \dots \\ y'_{m-1} &= y_m, \\ y'_m &= f(t, y_1, y_2, \dots, y_m). \end{aligned} \quad (14.128)$$

Начальные условия (14.127) в новых обозначениях принимают вид

$$y_1(t_0) = y_{10}, \quad y_2(t_0) = y_{20}, \dots, \quad y_m(t_0) = y_{m0}. \quad (14.129)$$

Пример 14.21. Задача Коши для дифференциального уравнения второго порядка $y'' = -25y + \sin t$, $y(0) = 0$, $y'(0) = 1$ введением новых

искомых функций $y_1(t) = y(t)$, $y_2(t) = y'(t)$ сводится к эквивалентной задаче Коши для системы дифференциальных уравнений первого порядка

$$y'_1 = y_2, \quad y'_2 = -25y_1 + \sin t, \quad y_1(0) = 0, \quad y_2(0) = 1. \quad \blacktriangle$$

Для решения задачи Коши (4.126), (4.127), приведенной к виду (14.128), (14.129), можно воспользоваться известными методами или даже готовыми программами. Часто именно так и поступают. Следует все же иметь в виду, что вычисления можно организовать и так, что сведение уравнения (14.126) к системе (14.128) не потребуется. Например, для решения дифференциального уравнения второго порядка $y'' = f(t, y)$ используется ряд специальных методов [109]. Одним из наиболее популярных среди них является *метод Нумерова* четвертого порядка точности:

$$\frac{y_{n+1} - 2y_n + y_{n-1}}{h^2} = \frac{1}{12} (f(t_{n-1}, y_{n-1}) + 10f(t_n, y_n) + f(t_{n+1}, y_{n+1})).$$

§ 14.11. Жесткие задачи

1. Понятие о жестких задачах. В последние годы при решении задачи Коши явными методами Рунге—Кутты и Адамса значительное число исследователей сталкивается с весьма неожиданным и неприятным явлением. Несмотря на медленное изменение искомых функций расчет приходится вести, казалось бы, с неоправданно мелким шагом h . Все попытки увеличить шаг и тем самым уменьшить время решения задачи приводят лишь к катастрофически большому росту погрешности. Обладающие таким свойством задачи получили название **жестких**. Сразу же подчеркнем, что жесткость является свойством задачи Коши (а не используемых численных методов).

Жесткие задачи встречаются в самых различных областях науки и техники. Традиционными источниками появления таких задач являются химическая кинетика, теория ядерных реакторов, теория автоматического управления, электротехника, электроника и т.д. Жесткие задачи возникают также при аппроксимации начально-краевых задач для уравнений в частных производных с помощью полудискретных методов (методов прямых).

Трудности, возникающие при численном решении жестких задач, продемонстрируем на таком примере.

Пример 14.22. Вычислим значения приближенного решения задачи Коши

$$y'(t) = -25y(t) + \cos t + 25\sin t, \quad y(0) = 1, \quad (14.130)$$

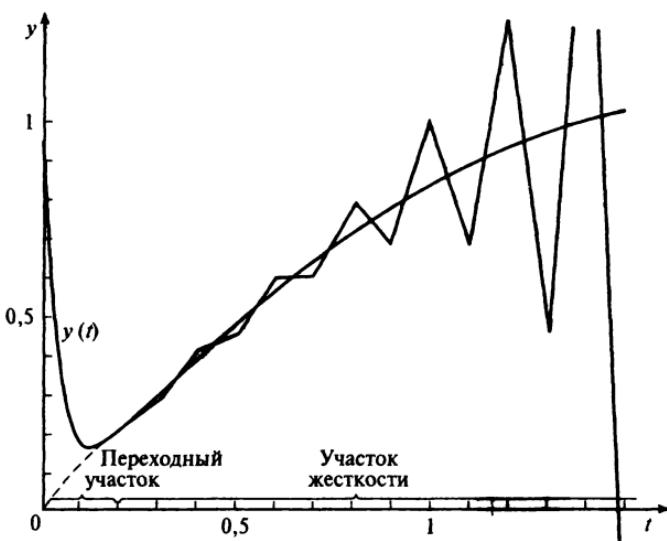


Рис. 14.16

использовав метод Эйлера. Решением этой задачи является функция $y(t) = \sin t + e^{-25t}$. Как нетрудно видеть (рис. 14.16), на начальном переходном участке $0 \leq t \leq 0.2$ решение быстро меняется. Однако уже через небольшой интервал времени переходная часть решения e^{-25t} практически исчезает и решение становится медленно меняющимся.

Естественно, что приближенное решение на переходном участке приходится вычислять, используя достаточно мелкий шаг. Однако при $t \geq 0.2$, когда переходная часть решения, казалось бы, уже практически отсутствует, возникает желание перейти к вычислению со сравнительно крупным шагом.

Предположим, что при $t = 0.2$ найдено значение решения $y(0.2) \approx 0.205$ с точностью $\epsilon = 10^{-3}$. Возьмем шаг $h = 0.1$ и будем вычислять решение при $t \geq 0.2$, используя метод Эйлера:

$$y_{n+1} = y_n + h(-25y_n + \cos t_n + 25\sin t_n).$$

Полученные значения приближений и точные значения решения приведены в первых трех столбцах табл. 14.6. Соответствующая ломаная Эйлера изображена на рис. 14.16. Как нетрудно видеть, метод ведет себя неустойчивым образом и оказывается непригодным для решения рассматриваемой задачи при значении шага $h = 0.1$. На конкретном примере мы убедились в том, что условие абсолютной устойчивости (14.96) нарушать нельзя. В рассматриваемом случае $\lambda = -25$ и это условие рав-

носильно требованию $h \leq 0.08$. Взяв удовлетворяющий этому условию шаг $h = 0.025$, получим вполне приемлемое решение (табл. 14.6).

Попробуем теперь воспользоваться для решения задачи (14.130) неявным методом Эйлера:

$$y_{n+1} = y_n + h(-25y_{n+1} + \cos t_{n+1} + 25\sin t_{n+1}).$$

Напомним, что он A -устойчив и, следовательно, абсолютно устойчив для всех значений h . Используя для нахождения решения формулу

$$y_{n+1} = \frac{y_n + h(\cos t_{n+1} + 25\sin t_{n+1})}{1 + 25h},$$

при $h = 0.1$ получаем значения, совпадающие с соответствующими значениями решения с точностью $2 \cdot 10^{-3}$ (табл. 14.6). Если же нас устраивает точность $\varepsilon = 6 \cdot 10^{-3}$, то шаг h можно утроить. Вычисленные с шагом $h = 0.3$ значения также указаны в табл. 14.6. ▲

Приведенная в примере ситуация типична для жестких задач. При использовании классических явных методов наличие в решении быстро меняющейся жесткой компоненты даже на том участке, где ее значение пренебрежимо мало, заставляет выбирать шаг h из условия абсолютной устойчивости. Для жестких задач это ограничение приводит к неприемлемо малому значению шага h , поэтому численное решение

Таблица 14.6

t	Точное решение	Метод Эйлера		Неявный метод Эйлера	
		$h = 0.1$	$h = 0.025$	$h = 0.1$	$h = 0.3$
0.2	0.205	0.205	0.205	0.205	0.205
0.3	0.296	0.287	0.296	0.304	—
0.4	0.389	0.404	0.390	0.391	—
0.5	0.479	0.460	0.480	0.479	0.478
0.6	0.565	0.597	0.565	0.564	—
0.7	0.644	0.600	0.644	0.643	—
0.8	0.717	0.787	0.718	0.716	0.714
0.9	0.783	0.683	0.784	0.782	—
1.0	0.841	0.996	0.842	0.840	—
1.1	0.891	0.664	0.892	0.890	0.886
1.2	0.932	1.277	0.932	0.930	—
1.3	0.964	0.451	0.964	0.962	—
1.4	0.985	1.759	0.986	0.984	0.980
1.5	0.997	-0.158	0.996	0.996	—

таких задач требует применения специальных неявных методов. Простейшим из них является неявный метод Эйлера. На рис. 14.17 и 14.18 проиллюстрировано принципиальное отличие результатов вычислений, осуществляемых с помощью явного и неявного методов Эйлера.

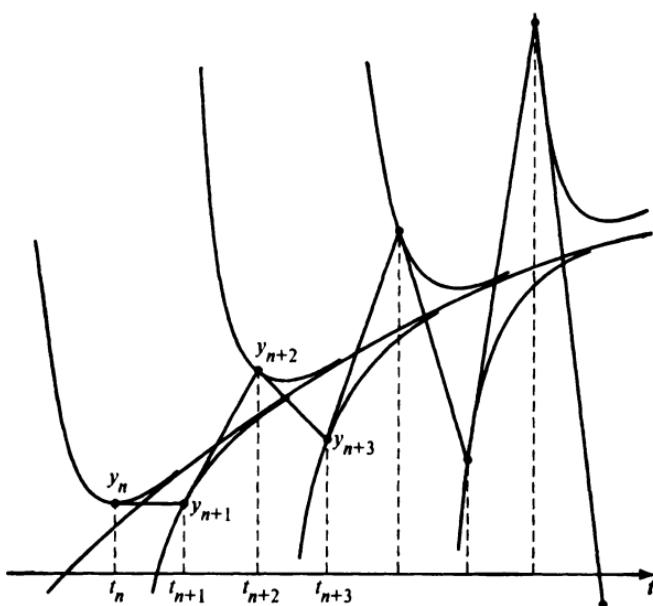


Рис. 14.17

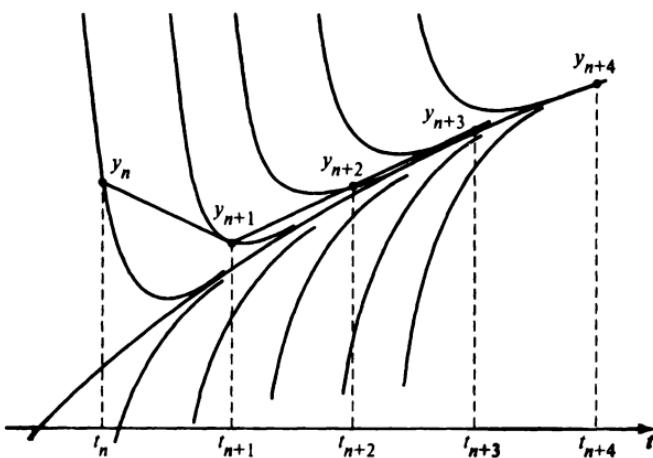


Рис. 14.18

2. Жесткие задачи для систем дифференциальных уравнений. Рассмотрим задачу Коши для системы линейных дифференциальных уравнений с постоянными коэффициентами

$$\mathbf{y}'(t) = \mathbf{A}\mathbf{y}(t). \quad (14.131)$$

Предположим, что \mathbf{A} — матрица простой структуры и все собственные числа этой матрицы имеют отрицательные вещественные части ($\operatorname{Re} \lambda_m \leq \operatorname{Re} \lambda_{m-1} \leq \dots \leq \operatorname{Re} \lambda_1 < 0$). Как отмечалось в § 14.10, в этом случае решение задачи асимптотически устойчиво и представляется в виде

$$\mathbf{y}(t) = c_1 e^{\lambda_1(t-t_0)} \mathbf{e}_1 + c_2 e^{\lambda_2(t-t_0)} \mathbf{e}_2 + \dots + c_m e^{\lambda_m(t-t_0)} \mathbf{e}_m.$$

Если среди собственных чисел λ_i имеются числа с сильно различающимися значениями вещественных частей, то возникает проблема, связанная с наличием в решении у компонент, имеющих существенно различные временные постоянные $\tau_i = 1/|\operatorname{Re} \lambda_i|$. Через довольно короткий интервал времени поведение решения будет определяться наиболее слабо меняющейся (*медленной*) компонентой решения. Так, если $\operatorname{Re} \lambda_m \leq \dots \leq \operatorname{Re} \lambda_2 < \operatorname{Re} \lambda_1 < 0$, то при $t - t_0 \gg \tau_2 = 1/|\operatorname{Re} \lambda_2|$ справедливо приближенное равенство $\mathbf{y}(t) \approx c_1 e^{\lambda_1(t-t_0)} \mathbf{e}_1$. В то же время применяемый для решения задачи Коши метод должен обладать свойствами устойчивости, позволяющими подавлять наиболее быстро меняющуюся (*жесткую*) компоненту $c_m e^{\lambda_m(t-t_0)} \mathbf{e}_m$ погрешности

$$\mathbf{e}(t) = \alpha_1 e^{\lambda_1(t-t_0)} \mathbf{e}_1 + \alpha_2 e^{\lambda_2(t-t_0)} \mathbf{e}_2 + \dots + \alpha_m e^{\lambda_m(t-t_0)} \mathbf{e}_m.$$

Отметим, что для медленной компоненты решения временной постоянной является величина $(\min_{1 \leq k \leq n} |\operatorname{Re} \lambda_k|)^{-1}$, а для жесткой компоненты погрешности — величина $(\max_{1 \leq k \leq n} |\operatorname{Re} \lambda_k|)^{-1}$. Их отношение и определяет степень жесткости задачи.

Приведем применительно к системе (14.131) одно из определений жесткости. Пусть $\operatorname{Re} \lambda_k < 0$ для всех $k = 1, \dots, m$. Определим число жесткости системы (14.131) с помощью формулы

$$s = \frac{\max_{1 \leq k \leq m} |\operatorname{Re} \lambda_k|}{\min_{1 \leq k \leq m} |\operatorname{Re} \lambda_k|}. \quad (14.132)$$

Систему уравнений (14.129) назовем *жесткой*, если для нее $s \gg 1$.

Пример 14.23. Рассмотрим систему

$$\begin{aligned}y' &= -7051y - 2499z, \\z' &= -7497y - 2503z, \\y(0) &= 0, \quad z(0) = 4.\end{aligned}$$

Собственные значения матрицы коэффициентов

$$A = \begin{bmatrix} -7501 & -2499 \\ -7497 & -2503 \end{bmatrix}$$

таковы: $\lambda_1 = -4$, $\lambda_2 = -10\,000$. Здесь число жесткости $s = |\lambda_2|/|\lambda_1| = 2500$ много больше единицы и поэтому систему можно квалифицировать как жесткую.

Общее решение системы имеет вид

$$y(t) = c_1 e^{-4t} + c_2 e^{-10000t}, \quad z(t) = -3c_1 e^{-4t} + c_2 e^{-10000t}.$$

Начальным значениям $y(0) = 0$, $z(0) = 4$ отвечает решение

$$y(t) = -e^{-4t} + e^{-10000t}, \quad z(t) = 3e^{-4t} + e^{-10000t}.$$

Жесткая компонента решения здесь быстро затухает и через очень небольшой интервал времени решение будет практически совпадать с $\tilde{y}(t) = -e^{-4t}$, $\tilde{z}(t) = 3e^{-4t}$.

Решение этой задачи с помощью явных методов очень неэффективно. Например, для метода Рунге—Кутты четвертого порядка точности условие устойчивости $h \leq \min_k \frac{2}{|\operatorname{Re} \lambda_k|}$ в данном случае приводит к необходимости использования очень мелкого шага $h \leq 2.8 \cdot 10^{-4}$.

Замечание. В приведенном выше наиболее простом определении жесткости требовалось, чтобы $\operatorname{Re} \lambda_k < 0$ для всех k . Более общие современные определения жесткости [9, 33, 108] не исключают наличия собственных чисел λ_i с положительными вещественными частями при условии, что для них $\operatorname{Re} \lambda_i \ll \max_{1 \leq k \leq m} |\operatorname{Re} \lambda_k|$.

Когда матрица A зависит от t (т.е. решается система линейных уравнений с переменными коэффициентами), число жесткости также зависит от t и определяется по формуле (14.132), в которой уже $\lambda_k = \lambda_k(t)$. Так как $s = s(t)$, то система $y' = A(t)y$ может оказаться жесткой на одном интервале времени t и нежесткой — на другом.

В настоящее время нет общепринятого математически корректного определения жесткости задачи Коши для системы нелинейных уравнений

$$y' = f(t, y). \quad (14.133)$$

Чаще всего эту задачу называют жесткой в окрестности точки $(\tilde{t}, y(\tilde{t}))$, если жесткой является соответствующая линеаризованная система

$$y' = Ay(t), \quad A = f'_y(\tilde{t}, y(\tilde{t})). \quad (14.134)$$

В подтверждение правомерности такого определения жесткости можно сослаться на то, что (см. § 14.10) погрешность $\varepsilon(t)$ решения нелинейной системы удовлетворяет приближенному равенству $\varepsilon'(t) \approx A\varepsilon(t)$ и поэтому в малой окрестности точки $(\tilde{t}, y(\tilde{t}))$ погрешности решений линеаризованной системы (14.134) и системы (14.133) должны вести себя примерно одинаковым образом.

К счастью, для того чтобы распознать жесткую задачу на практике, часто совсем не обязательно проводить математически строгое ее исследование. Если система уравнений (14.133) правильно моделирует реальное физическое явление, включающее процессы с существенно различными временными постоянными, то соответствующая задача Коши должна быть жесткой. Как правило, исследователь проявляет интерес к изучению поведения медленно меняющихся характеристик процесса в течение длительного времени. Наличие же быстро меняющихся физических компонент при использовании классических явных методов решения задачи Коши заставляет его выбирать шаг h порядка наименьшей из временных постоянных, что делает процесс численного решения чрезвычайно дорогостоящим и неэффективным. Существующие в настоящее время методы решения жестких задач позволяют использовать шаг h порядка наибольшей из временных постоянных, подчиняя его выбор только требованию точности.

Замечание 1. В последнее время задачу не принято квалифицировать как жесткую на переходном участке. Если исследователь проявляет интерес к изучению переходного режима, то для нахождения решения на переходном участке могут оказаться вполне приемлемыми и явные методы Рунге—Кутты и Адамса.

Замечание 2. Обычно при определении жесткости делается явное или неявное предположение о том, что среди собственных чисел матрицы A отсутствуют такие, для которых $|\operatorname{Im} \lambda_i| \gg 1$. Это означает, что предполагается отсутствие быстрых осцилляций в компонентах погрешности. Если же такие осцилляции возможны, то необходимо использовать специальные методы подавления соответствующих компонент погрешности.

3. Понятие о методах решения жестких задач. Для решения жестких задач было бы желательно использовать A -устойчивые методы, так как они не накладывают никаких ограничений на шаг h . Однако оказывается, что класс таких методов весьма узок. Например, среди явных линейных многошаговых методов нет A -устойчивых. Доказано также, что среди неявных линейных многошаговых методов нет A -устойчивых методов, имеющих порядок точности выше второго.

Многие из возникающих на практике жестких задач таковы, что для них собственные значения матрицы Якоби удовлетворяют неравенству $|\arg(-\lambda_i)| < \alpha$ ($i = 1, 2, \dots, m$), где $\alpha > 0$ — некоторое число. В частности, если все собственные значения вещественны и отрицательны, то указанное неравенство выполняется для любого сколь угодно малого $\alpha > 0$.

Для таких задач требование A -устойчивости методов является чрезмерным и его можно заменить менее ограничительным требованием наличия у метода $A(\alpha)$ -устойчивости. Численный метод решения задачи Коши называют $A(\alpha)$ -устойчивым, если область его абсолютной устойчивости включает угол $|\arg(-z)| < \alpha$ (рис. 14.19). В частности, при $\alpha = \pi/2$ определение $A(\alpha)$ -устойчивости совпадает с определением A -устойчивости.

Известно, что среди явных линейных многошаговых методов нет $A(\alpha)$ -устойчивых ни при каком $\alpha > 0$. Однако среди неявных линейных многошаговых методов имеются $A(\alpha)$ -устойчивые методы высокого порядка точности. Важный класс таких методов (*формул дифференцирования назад*) относится к так называемым *чисто неявным методам*.

$$\frac{1}{h} \sum_{j=0}^n \alpha_j y_{n+1-j} = f(t_{n+1}, y_{n+1}).$$

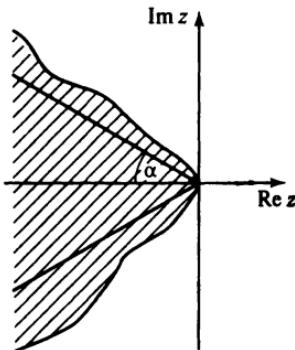


Рис. 14.19

Чисто неявные методы получаются в результате замены в системе дифференциальных уравнений (14.133) при $t = t_{n+1}$ производной $y'(t)$ ее разностной аппроксимацией, использующей значения функции в точках $t_{n+1}, t_n, \dots, t_{n+1-k}$. Если для этого используется односторонняя разностная производная (см. § 12.2), то получается *формула дифференцирования назад*.

Приведем формулы дифференцирования назад при $k = 1, 2, 3, 4$, имеющие k -й порядок точности

$$\frac{y_{n+1} - y_n}{h} = f(t_{n+1}, y_{n+1}), \quad k = 1 \quad (14.135)$$

(это неявный метод Эйлера);

$$\frac{3y_{n+1} - 4y_n + y_{n-1}}{2h} = f(t_{n+1}, y_{n+1}), \quad k = 2; \quad (14.136)$$

$$\frac{11y_{n+1} - 18y_n + 9y_{n-1} - 2y_{n-2}}{6h} = f(t_{n+1}, y_{n+1}), \quad k = 3; \quad (14.137)$$

$$\frac{25y_{n+1} - 48y_n + 36y_{n-1} - 16y_{n-2} + 3y_{n-3}}{12h} = f(t_{n+1}, y_{n+1}), \quad k = 4. \quad (14.138)$$

Напомним, что метод (14.135) является *A*-устойчивым. Как видно из рис. 14.20, *a* и *б*, для формул (14.136) и (14.137) области их абсолютной устойчивости почти целиком содержат левую полуплоскость $\operatorname{Re} z < 0$, поэтому свойства устойчивости¹ этих методов вполне достаточны для решения большинства жестких задач.

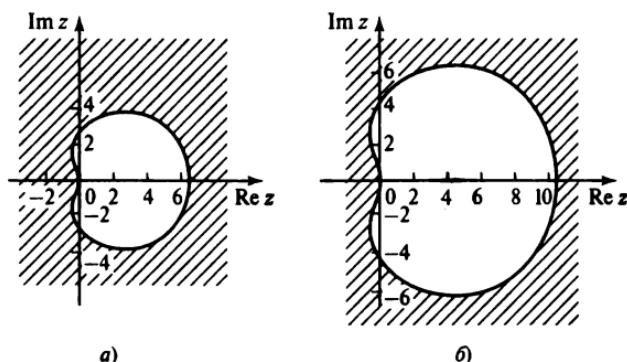


Рис. 14.20

¹ При $1 \leq k \leq 6$ формулы дифференцирования назад обладают так называемой *жесткой устойчивостью* [82, 91].

Весьма популярный и широко используемый при решении жестких задач *алгоритм Гира*¹ основан на использовании формул дифференцирования назад порядка точности $k = 1, 2, \dots, 6$ и представляет собой метод с автоматическим выбором шага интегрирования и порядка метода. Один из первых вариантов алгоритма реализован в фортранной программе DIFSUB.

§ 14.12. Дополнительные замечания

1. К настоящему времени разработано большое число различных численных методов решения задачи Коши и работа в этом направлении ведется очень активно. Тем не менее наиболее популярными остаются классические методы Рунге—Кутты и Адамса, а также их современные модификации. Каждый из этих двух классов методов имеет определенные достоинства и недостатки, некоторые из них уже обсуждались выше. Не имея перед собой конкретной задачи, вряд ли можно дать рекомендации в пользу того или иного метода, тем более что до сих пор в этом вопросе нет достаточной ясности. Однако ориентируясь на серьезное обсуждение оценки качества методов, приведенное в [91], можно попытаться грубо описать ситуации, в которых они обладают большей эффективностью. При этом следует иметь в виду, что одним из основных показателей эффективности метода является количество вычислений правых частей дифференциальных уравнений, которое требуется для достижения заданной точности решения.

Предположим, что решаемая задача Коши не является жесткой. Предположим также, что вычисление правых частей дифференциальных уравнений не является слишком трудоемкой операцией. Тогда целесообразно применение методов Рунге—Кутты с автоматическим выбором шага, наиболее эффективным среди которых для широкого класса задач являются методы Рунге—Кутты—Фельберга и Дормана—Принса пятого порядка точности. Если же к точности решения не предъявляются слишком высокие требования, то хороший результат следует ожидать и от применения классического метода Рунге—Кутты четвертого порядка точности.

В том случае, когда вычисления правых частей трудоемки, имеет смысл предпочесть использование качественной программы, реализующей метод Адамса с автоматическим выбором шага и порядка метода.

¹ Чарльз Вильям (Билл) Гир (1935) — английский математик. В 1971 г. предложил новый класс методов решения жестких задач, принесший ему заслуженную известность. В 1986 г. Гир был избран президентом Общества промышленной и прикладной математики (Society for Industrial and Applied Mathematics – SIAM).

По-видимому, именно эти методы в будущем станут наиболее употребительными.

2. После того как приближенные значения y_n решения задачи Коши в узлах t_n ($n = 0, 1, \dots, N$) определены, для вычисления значений решения $y(t)$ в промежуточных точках можно использовать интерполяцию. В связи с этим полезно отметить, что наряду со значениями вектор-функции y_n фактически оказываются вычисленными значения производной $y'_n = f(t_n, y_n)$. Поэтому в данном случае для интерполяции естественно использование кубического интерполяционного многочлена Эрмита или локального кубического сплайна (см. гл. 11). Для методов первого или второго порядка точности вполне удовлетворительный результат дает использование линейной интерполяции.

3. В последние десятилетия все большую популярность получают так называемые *многозначные методы* (*методы Нордсика*).

Как уже отмечалось выше, при реализации многошаговых методов возникают две основные проблемы. Первая из них связана с нестандартным определением начальных значений. Вторая проблема возникает при попытке изменить шаг интегрирования. Многозначные методы свободны от этих проблем, так как они являются одношаговыми. Тем не менее их можно рассматривать как особую форму реализации многошаговых методов.

В $(k+1)$ -значных методах в каждом из узлов t_n вычисляется вектор Y_n (который называют *вектором Нордсика*), содержащий не только приближения к значениям $y(t_n)$, но и к значениям дифференциалов

$$dy(t_n, h) = y'(t_n)h, \quad d^2y(t_n, h) = y''(t_n) \frac{h^2}{2}, \quad \dots, \quad d^k y(t_n, h) = y^{(k)}(t_n) \frac{h^k}{k!}.$$

Таким образом, $(k+1)$ -значный метод по значению вектора

$$Y_n \approx \left(y(t_n), y'(t_n)h, y''(t_n) \frac{h^2}{2}, \dots, y^{(k)}(t_n) \frac{h^k}{k!} \right)^T$$

дает приближение к вектору

$$Y_{n+1} \approx \left(y(t_{n+1}), y'(t_{n+1})h, y''(t_{n+1}) \frac{h^2}{2}, \dots, y^{(k)}(t_{n+1}) \frac{h^k}{k!} \right)^T$$

Подробнее с многозначными методами (методами Нордсика) можно познакомиться в [109, 15*].

4. Поиск эффективных методов решения жестких задач еще находится в начальной стадии. Тем не менее разработан ряд популярных алгоритмов (среди которых наиболее известен алгоритм Гира) и соз-

дано значительное число качественных программ. Вопрос о наиболее эффективном методе решения жестких задач остается открытым и какие-либо рекомендации здесь преждевременны.

В последнее время выяснилась достаточная перспективность применения для решения жестких задач некоторых неявных методов Рунге—Кутты [33, 109, 15*]. И, хотя классические явные методы Рунге—Кутты для этой цели совершенно непригодны, разработаны специальные явные методы. Рунге—Кутты с областью абсолютной устойчивости, расширенной вдоль отрицательной части вещественной оси. Обсуждение этих методов, наиболее известным среди которых является метод, предложенный в 1989 г. В.И. Лебедевым¹, можно найти в [108].

4. Дополнительную информацию о методах решения задачи Коши (и полезное обсуждение жестких задач) можно найти, например в учебниках [9, 16, 51, 73, 84, 88, 106]. Настоятельно советуем обратить внимание на весьма содержательные книги [91, 109, 15*]. Вторая из них содержит систематическое и доступное широкому кругу читателей изложение численных методов решения нежестких задач. Методам решения жестких задач специально посвящены монографии [33, 82, 108]; последние две из них содержат современный взгляд на эту проблему.

¹ Вячеслав Иванович Лебедев (1930) — российский математик.

РЕШЕНИЕ ДВУХТОЧЕЧНЫХ КРАЕВЫХ ЗАДАЧ

Двухточечная краевая задача — это задача отыскания решения обыкновенного дифференциального уравнения или системы обыкновенных дифференциальных уравнений на отрезке $a \leq x \leq b$, в которой дополнительные условия на решение налагаются в двух точках a и b — «краях» отрезка (отсюда — и название задачи).

Решить краевую задачу, вообще говоря, значительно труднее, чем задачу Коши и для этого используются разнообразные подходы. Наиболее распространены различные методы дискретизации, позволяющие заменить исходную задачу некоторым ее дискретным аналогом. Получающаяся дискретная краевая задача представляет собой систему уравнений (возможно, нелинейных) с конечным числом неизвестных и может быть решена на компьютере с помощью специальных прямых или итерационных методов. Одним из простейших и весьма популярных подходов к дискретизации является использование метода конечных разностей. В § 15.2 и 15.3 рассматриваются некоторые из основных моментов применения этого метода.

В § 15.4 дается представление о другом подходе к дискретизации краевых задач. В нем описываются проекционные методы Ритца и Галеркина и обсуждается один из их современных вариантов, имеющий большое практическое значение, — метод конечных элементов.

В заключение главы рассматривается метод пристрелки.

§ 15.1. Краевые задачи для одномерного стационарного уравнения теплопроводности

1. Дифференциальное уравнение и краевые условия. Рассмотрим дифференциальное уравнение второго порядка

$$-\frac{d}{dx} \left(k(x) \frac{du}{dx} \right) + q(x)u = f(x), \quad a \leq x \leq b. \quad (15.1)$$

Оно называется *одномерным стационарным уравнением теплопроводности* и возникает при математическом моделировании многих важных процессов. Например, это уравнение описывает установившееся распределение температуры $u(x)$ в теплопроводящем стержне длины $l = b - a$. В этом случае $k(x)$ — коэффициент теплопроводности; $w(x) =$

$= -k(x) \frac{du}{dx}$ — плотность потока тепла; $q(x)$ — коэффициент теплоотдачи (q_u — мощность стоков тепла, пропорциональная температуре u); $f(x)$ — плотность источников тепла (при $f \leq 0$ — плотность стоков тепла).

Уравнение (15.1) описывает также установившееся распределение плотности нейтронов в реакторе, характеристики которого зависят от одной пространственной переменной x . В такой трактовке $u(x)$ — это полный поток нейтронов, $k(x)$ — коэффициент диффузии, $q(x)$ — сечение поглощения, $f(x)$ — плотность источников нейтронов. То же уравнение описывает и стационарные процессы диффузии газов (растворов) в пористых средах, $u(x)$ рассматривается тогда как концентрация диффундирующего вещества. Поэтому уравнение (15.1) часто называют *одномерным уравнением диффузии*. Рассматриваемое уравнение имеет приложения и в других областях техники и естествознания (деформации струн и стержней, распространение электромагнитных волн и т.д.).

Далее будем считать функции $k(x)$, $q(x)$, $f(x)$ заданными и предполагать, что выполнены неравенства

$$k(x) \geq k_0 > 0, \quad q(x) \geq 0. \quad (15.2)$$

Так как уравнение (15.1) является дифференциальным уравнением второго порядка, то для того чтобы однозначно найти функцию $u(x)$ — распределение температуры в стержне, необходимо задать два дополнительных условия. Простейшая постановка краевых условий такова:

$$u(a) = u_a, \quad u(b) = u_b.$$

Краевые условия такого типа принято называть *краевыми условиями первого рода*. Физическая интерпретация этих краевых условий состоит в том, что в рассматриваемой задаче на торцах стержня поддерживаются фиксированные значения температуры u_a и u_b .

Возможны и другие постановки краевых условий. Так, если известна плотность потока тепла через левый торец стержня, то условие $u(a) = u_a$ можно заменить *краевым условием второго рода*:

$$-k(a)u'(a) = w_a.$$

Аналогичное условие для правого торца имеет вид

$$-k(b)u'(b) = w_b.$$

Основное внимание в этой главе будет уделено краевой задаче

$$L[u](x) = f(x), \quad a < x < b; \quad (15.3)$$

$$u(a) = u_a, \quad u(b) = u_b. \quad (15.4)$$

Здесь L — дифференциальный оператор, определяемый следующим образом:

$$L[u](x) = -(k(x)u'(x))' + q(x)u(x).$$

2. Разрешимость краевой задачи. Будем считать, что коэффициенты q и f непрерывны на отрезке $[a, b]$, коэффициент k непрерывно дифференцируем на $[a, b]$ и выполнены условия (15.2).

Назовем дважды непрерывно дифференцируемую на отрезке $[a, b]$ функцию $u(x)$ *решением (классическим решением) краевой задачи* (15.3), (15.4), если $u(x)$ является решением дифференциального уравнения (15.3) и удовлетворяет краевым условиям (15.4).

Приведем без доказательства известные из теории дифференциальных уравнений результаты о разрешимости рассматриваемой краевой задачи и о гладкости ее решения.

Теорема 15.1. *Решение краевой задачи (15.3), (15.4) существует и единственно.* ■

Теорема 15.2. *Пусть коэффициенты q и f являются m раз, а коэффициент $k - m + 1$ раз непрерывно дифференцируемыми на отрезке $[a, b]$ функциями. Тогда решение и краевой задачи (15.3), (15.4) является $m + 2$ раза непрерывно дифференцируемой на отрезке $[a, b]$ функцией.* ■

3. Принцип максимума. Важным свойством уравнения (15.3) является наличие так называемого *принципа максимума*. Приведем один из вариантов его формулировки.

Теорема 15.3. *Пусть $u(x)$ — решение задачи (15.3), (15.4). Тогда если $f(x) \leq 0$, $u_a \leq 0$, $u_b \leq 0$, то $u(x) \leq 0$.* ■

Теорема 15.3 имеет простой физический смысл. Если отсутствуют источники тепла и температура торцов стержня неположительна, то ни в одной из внутренних точек стержня температура не может стать положительной.

Заметим, что произвольную дважды непрерывно дифференцируемую функцию $u(x)$ можно рассматривать как решение краевой задачи (15.3), (15.4), если специальным образом выбрать правую часть f и краевые значения u_a , u_b , а именно положить $f(x) = L[u](x)$, $u_a = u(a)$, $u_b = u(b)$. С учетом этого замечания, сформулируем теорему 15.3 иным образом.

Теорема 15.4. *Пусть $u(x)$ — дважды непрерывно дифференцируемая на отрезке $[a, b]$ функция, удовлетворяющая неравенствам $L[u] \leq 0$, $u(a) \leq 0$, $u(b) \leq 0$. Тогда $u(x) \leq 0$.* ■

Из теоремы 15.4 вытекает следующее утверждение.

Теорема 15.5 (теорема сравнения). Пусть $u(x), v(x)$ — дважды непрерывно дифференцируемые на отрезке $[a, b]$ функции, удовлетворяющие неравенствам $L[u] \leq L[v]$, $u(a) \leq v(a)$, $u(b) \leq v(b)$. Тогда $u(x) \leq v(x)$. ■

4. Априорная оценка и устойчивость решения. Используя теорему сравнения, можно вывести оценку максимума модуля решения $u(x)$ через данные краевой задачи.

Теорема 15.6. Справедлива следующая оценка решения краевой задачи (15.3), (15.4):

$$\max_{[a, b]} |u(x)| \leq \max \{ |u_a|, |u_b| \} + K \cdot \max_{[a, b]} |f(x)|. \quad (15.5)$$

Здесь $K = \frac{Rl}{4}$, $R = \int_a^b \frac{dx}{k(x)}$, $l = b - a$. ■

Замечание 1. Неравенства типа (15.5) принято называть *априорными оценками решения*.

Замечание 2. Если коэффициент $k(x)$ рассматривать как коэффициент теплопроводности, то $\frac{1}{k(x)}$ — это коэффициент теплосопротивления, а $R = \int_a^b \frac{dx}{k(x)}$ — это полное теплосопротивление стержня.

Замечание 3. При $k(x) = 1$ уравнение (15.3) принимает вид

$$-u''(x) + q(x)u(x) = f(x). \quad (15.6)$$

В этом случае $R = l$ и оценку (15.5) можно уточнить следующим образом:

$$\max_{[a, b]} |u(x)| \leq \max \{ |u_a|, |u_b| \} + \frac{l^2}{8} \max_{[a, b]} |f(x)|. \quad (15.7)$$

Рассмотрим теперь вопрос о влиянии погрешностей задания краевых значений u_a , u_b и правой части f на решение краевой задачи. Пусть $u(x)$ — решение краевой задачи (15.3), (15.4), а $u^*(x)$ — решение краевой задачи

$$L[u^*](x) = f^*(x), \quad a < x < b;$$

$$u^*(a) = u_a^*, \quad u^*(b) = u_b^*.$$

Здесь $f^*(x)$ — непрерывная функция, рассматриваемая как приближенно заданная (с погрешностью $\delta f(x) = f(x) - f^*(x)$) правая часть уравнения; u_a^* , u_b^* — приближенно заданные (с погрешностями $\varepsilon_a = u_a - u_a^*$, $\varepsilon_b = u_b - u_b^*$) краевые значения.

Теорема 15.7. Справедлива оценка

$$\max_{[a, b]} |u(x) - u^*(x)| \leq \max \{ |\varepsilon_a|, |\varepsilon_b| \} + K \cdot \max_{[a, b]} |\delta f(x)|, \quad (15.8)$$

где K — та же постоянная что и в неравенстве (15.5).

Доказательство. Для доказательства достаточно заметить, что функция $\varepsilon(x) = u(x) - u^*(x)$ является решением краевой задачи

$$L[\varepsilon](x) = \delta f(x), \quad a < x < b;$$

$$\varepsilon(a) = \varepsilon_a, \quad \varepsilon(b) = \varepsilon_b,$$

и воспользоваться для оценки величины $\max_{[a, b]} |\varepsilon(x)|$ теоремой 15.6. ■

Из оценки (15.8) видно, что когда значение K не очень велико, краевая задача (15.3), (15.4) хорошо обусловлена. Если же $K \gg 1$, то задача является плохо обусловленной. В этом случае погрешности порядка δ задания правой части уравнений может отвечать погрешность порядка $K\delta$ решения задачи. Ниже при рассмотрении численных методов решения краевой задачи будем предполагать, что она хорошо обусловлена.

Замечание. Если рассматривается устойчивость решения краевой задачи для уравнения (15.6), то в оценке (15.8) следует заменить K на $P/8$.

§ 15.2. Метод конечных разностей: основные понятия

Метод конечных разностей (или *метод сеток*) является одним из универсальных и широко используемых методов решения краевых задач. Его популярность во многом объясняется относительной простотой подхода к дискретизации дифференциальных уравнений. Суть метода состоит в следующем. Область непрерывного изменения аргумента заменяют конечным (дискретным) множеством точек (узлов), называемым *сеткой*. Вместо функций непрерывного аргумента рассматривают функции, определенные только в узлах сетки, — *сеточные функции*. Производные, которые входят в дифференциальное уравнение

и краевые условия, заменяют их разностными аналогами — линейными комбинациями значений сеточных функций в некоторых узлах сетки. В результате краевую задачу заменяют *дискретной краевой задачей (разностной схемой)*, представляющей собой систему конечного числа линейных или нелинейных уравнений. Решение разностной схемы (предполагается, что оно существует) принимают за приближенное решение краевой задачи.

Несмотря на кажущуюся простоту метода при его использовании приходится решать ряд проблем. Например, следует иметь в виду, что для одной краевой задачи можно построить большое число различных разностных схем, среди которых далеко не все пригодны для использования на практике.

В этом параграфе мы покажем, как применяется разностный метод для решения краевой задачи (15.3), (15.4), ограничившись для простоты изложения уравнением с постоянным коэффициентом $k(x) = 1$. В этом случае краевая задача принимает вид

$$-u''(x) + q(x)u(x) = f(x), \quad a < x < b; \quad (15.9)$$

$$u(x) = u_a, \quad u(b) = u_b. \quad (15.10)$$

1. Построение сетки и введение сеточных функций. Произведем дискретизацию области непрерывного изменения аргумента x , заменив отрезок $[a, b]$ сеткой $\bar{\omega}^h$ — конечным набором точек $a = x_0 < x_1 < \dots < x_N = b$ (рис. 15.1). Точки x_i называются *узлами сетки* $\bar{\omega}^h$.

Для упрощения изложения в этом параграфе будем считать сетку равномерной с шагом $h = (b - a)/N$. Тогда $x_i - x_{i-1} = h$ для всех $i = 1, 2, \dots, N$. Заметим, что при $N \rightarrow \infty$ шаг $h \rightarrow 0$ (сетка измельчается).

Сетка $\bar{\omega}^h$ естественным образом разбивается здесь на два подмножества: $\bar{\omega}^h = \omega^h \cup \gamma^h$. *Множество внутренних узлов* ω^h состоит из тех узлов x_i ($1 \leq i \leq N - 1$), которые лежат внутри интервала (a, b) . *Множество граничных узлов* γ^h состоит из двух узлов $x_0 = a$ и $x_N = b$, лежащих на границе отрезка $[a, b]$.

Далее будем вычислять решение краевой задачи не в произвольных точках отрезка $[a, b]$, а только в узлах сетки $\bar{\omega}^h$. Таким образом, иско-

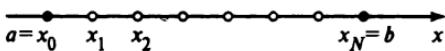


Рис. 15.1

мой окажется не функция u , а сеточная функция¹ u^h . Значения $u^h(x_i)$ этой функции в узлах x_i будем обозначать через u_i и рассматривать как приближения к значениям $u(x_i)$ решения задачи (15.9), (15.10). Введем также сеточные функции q^h и f^h , принимающие в узлах сетки $\bar{\omega}^h$ значения $q_i = q(x_i)$ и $f_i = f(x_i)$.

2. Построение разностной схемы. Напомним (см. гл. 12), что производную $u''(x)$ можно аппроксимировать второй разностной производной:

$$u''(x) \approx \frac{u(x-h) - 2u(x) + u(x+h)}{h^2} \quad (15.11)$$

с погрешностью $r_h(x) = u^{(4)}(\xi) \frac{h^2}{12}$, где $\xi \in [x-h, x+h]$. Используя формулу (15.11), заменим теперь в каждом из внутренних узлов x_i ($1 \leq i \leq N-1$) дифференциальное уравнение (15.9) приближенным равенством

$$-\frac{u(x_{i-1}) - 2u(x_i) + u(x_{i+1})}{h^2} + q_i u(x_i) \approx f_i, \quad 1 \leq i \leq N-1, \quad (15.12)$$

связывающим неизвестные значения решения в трех последовательных узлах сетки.

Потребуем теперь, чтобы значения искомой сеточной функции u^h удовлетворяли во всех внутренних узлах сетки уравнениям (15.12), в которых знак приближенного равенства заменен на знак равенства:

$$-\frac{u^h(x_{i-1}) - 2u^h(x_i) + u^h(x_{i+1})}{h^2} + q_i u^h(x_i) = f_i, \quad 1 \leq i \leq N-1. \quad (15.13)$$

В результате дифференциальное уравнение (15.9) оказалось аппроксимированным его дискретным аналогом — *разностным уравнением* (15.13).

Естественно потребовать, чтобы в граничных узлах сеточная функция u^h удовлетворяла равенствам

$$u^h(x_0) = u_a, \quad u^h(x_N) = u_b. \quad (15.14)$$

Таким образом, мы пришли к системе линейных алгебраических уравнений (15.13), (15.14), в которой число уравнений совпадает с числом неизвестных $u_i = u^h(x_i)$ ($i = 0, 1, \dots, N$) и равно $N + 1$. Решив эту систему (которую мы будем называть *системой сеточных уравнений*), можно найти сеточную функцию u^h .

¹ Сеточные функции использовались нами и в предыдущей главе.

Введем линейный разностный оператор L^h с помощью равенства

$$L^h[u^h](x_i) = -\frac{u^h(x_{i-1}) - 2u^h(x_i) + u^h(x_{i+1})}{h^2} + q_i u^h(x_i), \quad x_i \in \omega^h$$

и запишем систему сеточных уравнений (15.13), (15.14) в следующем виде:

$$L^h[u^h](x) = f^h(x), \quad x \in \omega^h; \quad (15.15)$$

$$u^h(x_0) = u_a, \quad u^h(x_N) = u_b. \quad (15.16)$$

Дискретную задачу (15.15), (15.16), зависящую от параметра h , принято называть *разностной схемой для краевой задачи* (15.9), (15.10).

3. Вычисление решения разностной схемы с помощью метода прогонки.

Приведем систему сеточных уравнений к виду

$$-u_{i-1} + (2 + h^2 q_i) u_i - u_{i+1} = h^2 f_i, \quad 1 \leq i \leq N-1; \quad (15.17)$$

$$u_0 = u_a, \quad u_N = u_b. \quad (15.18)$$

Как нетрудно видеть, эта система есть частный случай системы линейных алгебраических уравнений вида

$$\begin{aligned} b_0 u_0 + c_0 u_1 &= d_0, \\ a_i u_{i-1} + b_i u_i + c_i u_{i+1} &= d_i, \quad 1 \leq i \leq N-1; \\ a_N u_{N-1} + b_N u_N &= d_N, \end{aligned}$$

матрица которой трехдиагональна. Здесь

$$\begin{aligned} a_i &= -1, \quad b_i = 2 + h^2 q_i, \quad c_i = -1, \quad d_i = h^2 f_i, \quad 1 \leq i < N; \\ b_0 &= 1, \quad c_0 = 0, \quad d_0 = u_a, \quad a_N = 0, \quad b_N = 1, \quad d_N = u_b. \end{aligned} \quad (15.19)$$

Напомним, что эффективным методом решения таких систем является метод прогонки (см. § 5.9), вычисления которого состоят из двух этапов: прямого и обратного хода.

Прямой ход метода прогонки заключается в вычислении прогоночных коэффициентов α_i и β_i ($0 \leq i \leq N$). При $i = 0$ коэффициенты вычисляют по формулам

$$\alpha_0 = -c_0 / \gamma_0, \quad \beta_0 = d_0 / \gamma_0, \quad \gamma_0 = b_0,$$

а при $i = 1, 2, \dots, N-1$ — по рекуррентным формулам

$$\alpha_i = -c_i / \gamma_i, \quad \beta_i = (d_i - a_i \beta_{i-1}) / \gamma_i, \quad \gamma_i = b_i + a_i \alpha_{i-1}.$$

При $i = N$ прямой ход завершают вычислением значения

$$\beta_N = (d_N - a_N \beta_{N-1}) / \gamma_N, \quad \gamma_N = b_N + a_N \alpha_{N-1}.$$

Обратный ход метода прогонки дает значения неизвестных. Сначала полагают $u_N = \beta_N$, а затем значения остальных неизвестных находят по формуле

$$u_i = \alpha_i u_{i+1} + \beta_i.$$

Вычисления ведут в порядке убывания значений индекса i от $N - 1$ до 0.

Применительно к решению системы (15.17), (15.18) расчетные формулы метода прогонки упрощаются.

Прогоночные коэффициенты вычисляют по формулам

$$\begin{aligned}\alpha_0 &= 0, \quad \beta_0 = u_a, \\ \alpha_i &= 1/\gamma_i, \quad \beta_i = (h^2 f_i + \beta_{i-1})/\gamma_i, \quad \gamma_i = 2 + h^2 q_i - \alpha_{i-1},\end{aligned}\tag{15.20}$$

где $i = 1, 2, \dots, N - 1$.

Обратный ход дает значения неизвестных u_1, u_2, \dots, u_{N-1} (напомним, что значения $u_0 = u_a, u_N = u_b$ известны). Для этого производят вычисления по формулам

$$u_N = u_b,\tag{15.21}$$

$$u_i = \alpha_i u_{i+1} + \beta_i, \quad i = N - 1, N - 2, \dots, 1.\tag{15.22}$$

Замечание. Очевидно, что коэффициенты (15.19) удовлетворяют неравенствам

$$a_i \leq 0, \quad b_i > 0, \quad c_i < 0, \quad a_i + b_i + c_i \geq 0, \quad 1 \leq i < N.\tag{15.23}$$

Отсюда следует, что для системы (15.17), (15.18) выполнены условия диагонального преобладания: $|b_0| \geq |c_0|, |b_i| \geq |a_i| + |c_i| > |a_i|$ ($1 \leq i < N$), $|b_N| > |a_N|$. Поэтому в силу теоремы 5.4 вычисления по формулам (15.20) могут быть доведены до конца (ни один из знаменателей γ_i не обратится в нуль). Кроме того, обратная прогонка устойчива по входным данным.

4. Существование и единственность решения. Согласно последнему замечанию, систему сеточных уравнений (15.17), (15.18) с помощью эквивалентных преобразований можно привести к системе вида (15.21), (15.22), из которой однозначно находятся неизвестные. Таким образом, справедлива следующая теорема.

Теорема 15.8. Решение разностной схемы (15.15), (15.16) существует и единствено. ■

Как будет показано ниже, разностная схема (15.15), (15.16) обладает рядом свойств, аналогичных соответствующим свойствам краевой задачи (15.9), (15.10).

5. Принцип максимума. Как уже отмечалось ранее, важным свойством задачи (15.9), (15.10) является принцип максимума. Естественно потребовать, чтобы и для разностной схемы был справедлив дискретный аналог этого свойства. Более того, невыполнение для разностной схемы принципа максимума можно рассматривать как серьезный дефект, ставящий под сомнение возможность ее использования на практике.

Лемма 15.1 (принцип максимума для системы сеточных уравнений). Пусть сеточная функция u^h является решением системы сеточных уравнений

$$a_i u_{i-1} + b_i u_i + c_i u_{i+1} = d_i, \quad 1 \leq i < N; \quad (15.24)$$

$$u_0 = u_a, \quad u_N = u_b,$$

коэффициенты которой удовлетворяют условиям (15.23). Тогда если $u_a \leq 0$, $u_b \leq 0$ и $d_i \leq 0$ для всех $i = 1, 2, \dots, N - 1$, то $u^h \leq 0$.

Доказательство. Предположим, что неравенство $u_p \leq 0$ не выполнено. Так как по условию значения u^h в граничных узлах неотрицательны ($u_0 = u_a \leq 0$, $u_N = u_b \leq 0$), то максимальное значение функции u^h положительно и достигается во внутреннем узле сетки: $u_{\max} = \max_{1 \leq i < N} u_i > 0$.

Пусть j — максимальный среди индексов i , для которых $u_i = u_{\max}$. В силу такого выбора j справедливы неравенства $u_{j-1} \leq u_j = u_{\max}$, $u_{j+1} < u_j = u_{\max}$. Так как $a_j \leq 0$, $c_j < 0$, то $a_j u_j \leq a_j u_{j-1}$, $c_j u_j < c_j u_{j+1}$. Учитывая что $a_j + b_j + c_j \geq 0$, $d_j \leq 0$, из равенства (15.24), взятого при $i = j$, получаем следующую цепочку неравенств:

$$0 \leq (a_j + b_j + c_j)u_j < a_j u_{j-1} + b_j u_j + c_j u_{j+1} = d_j \leq 0.$$

Полученное противоречие ($0 < 0$) доказывает, что $u^h \leq 0$. ■

Теорема 15.9 (принцип максимума для разностной схемы). Пусть сеточная функция u^h является решением разностной схемы (15.15), (15.16). Тогда если $f^h \leq 0$, $u_a \leq 0$, $u_b \leq 0$, то $u^h \leq 0$.

Доказательство. Для доказательства достаточно заметить, что коэффициенты соответствующей системы сеточных уравнений (15.17), (15.18) удовлетворяют условиям (15.23), $d_i = h^2 f_i \leq 0$, и воспользоваться леммой 15.1. ■

Заметим, что произвольную сеточную функцию u^h можно считать решением разностной схемы (15.15), (15.16), если выбрать правую часть и граничные значения специальным образом, а именно положить $f^h = L^h[u^h]$, $u_a = u^h(x_0)$, $u_b = u^h(x_N)$. С учетом этого замечания сформулируем теорему 15.9 иным образом:

Теорема 15.10. *Пусть сеточная функция u^h удовлетворяет неравенствам $L^h[u^h] \leq 0$, $u^h(x_0) \leq 0$, $u^h(x_N) \leq 0$. Тогда $u^h \leq 0$. ■*

Из этой теоремы вытекает следующий важный результат.

Теорема 15.11 (теорема сравнения). *Пусть сеточные функции u^h и v^h удовлетворяют неравенствам $|L^h[u^h]| \leq L^h[v^h]$, $|u^h(x_0)| \leq v^h(x_0)$, $|u^h(x_N)| \leq v^h(x_N)$. Тогда $|u^h| \leq v^h$.*

Доказательство. Согласно условию, сеточная функция $y^h = u^h - v^h$ удовлетворяет неравенствам

$$L^h[y^h] = L^h[u^h] - L^h[v^h] \leq 0,$$

$$y^h(x_0) = u^h(x_0) - v^h(x_0) \leq 0, \quad y^h(x_N) = u^h(x_N) - v^h(x_N) \leq 0.$$

Поэтому в силу теоремы 15.10 $y^h \leq 0$, что эквивалентно неравенству $u^h \leq v^h$.

Аналогично, $z^h = -v^h - u^h$ удовлетворяет неравенствам

$$L^h[z^h] = -L^h[v^h] - L^h[u^h] \leq 0,$$

$$z^h(x_0) = -v^h(x_0) - u^h(x_0) \leq 0, \quad z^h(x_N) = -v^h(x_N) - u^h(x_N) \leq 0.$$

Следовательно, $z^h \leq 0$, что эквивалентно неравенству $-v^h \leq u^h$.

Итак, $-v^h \leq u^h \leq v^h$. ■

6. Априорная оценка решения. Оценим максимум модуля решения u^h разностной схемы через данные дискретной задачи (правую часть уравнения и краевые значения). Предварительно установим справедливость двух вспомогательных утверждений.

Лемма 15.2. *Для решения разностной схемы*

$$L^h[y^h] = 0, \quad x \in \omega^h; \tag{15.25}$$

$$y^h(x_0) = u_a, \quad y^h(x_N) = u_b \tag{15.26}$$

справедлива оценка

$$\max_{0 \leq i < N} |y_i| \leq \max \{ |u_a|, |u_b| \}. \tag{15.27}$$

Доказательство. Введем сеточную функцию $v^h \equiv M = \text{const}$, $M = \max\{|u_a|, |u_b|\}$. Заметим, что

$$|L^h[y^h]| = 0 \leq L^h[v^h] = q^h M,$$

$$|y^h(x_0)| = |u_a| \leq M = v^h(x_0), \quad |y^h(x_N)| = |u_b| \leq M = v^h(x_N).$$

Поэтому в силу теоремы сравнения $|y^h| \leq M$, что эквивалентно оценке (15.27). ■

Лемма 15.3. Для решения разностной схемы

$$L^h[z^h] = f^h, \quad x \in \omega^h; \quad (15.28)$$

$$z^h(x_0) = 0, \quad z^h(x_N) = 0 \quad (15.29)$$

справедлива оценка

$$\max_{0 \leq i \leq N} |z_i| \leq \frac{l^2}{8} \max_{0 < i < N} |f_i|. \quad (15.30)$$

Доказательство. Введем сеточную функцию

$$v^h(x_i) = \frac{A}{2} (x_i - a)(b - x_i),$$

где $A = \max_{0 < i < N} |f_i|$.

Заметим, что $v^h \geq 0$. Непосредственной проверкой нетрудно убедиться в том, что $L^h[v^h] = A + q^h v^h$. Таким образом,

$$|L^h[z^h]| = |f^h| \leq A \leq L^h[v^h],$$

$$|z^h(x_0)| = 0 \leq v^h(x_0), \quad |z^h(x_N)| = 0 \leq v^h(x_N)$$

и поэтому в силу теоремы сравнения $|z^h| \leq v^h$. Так как максимум квадратичной функции $\frac{A}{2} (x - a)(b - x)$ достигается при $x = \frac{a+b}{2}$ и равен $\frac{A}{8} l^2$,

то $v^h \leq \frac{A}{8} l^2$ и из неравенства $|z^h| \leq v^h$ следует оценка (15.30). ■

Сформулируем теперь основной результат этого пункта.

Теорема 15.12. Для решения разностной схемы (15.15), (15.16) справедлива априорная оценка

$$\max_{0 \leq i \leq N} |u_i| \leq \max\{|u_a|, |u_b|\} + \frac{l^2}{8} \max_{0 < i < N} |f_i| \quad (15.31)$$

Доказательство. Заметим, что сеточную функцию u^h можно представить в виде суммы: $u^h = y^h + z^h$, где y^h — решение разностной схемы (15.25), (15.26), а z^h — решение разностной схемы (15.28),

(15.29). Пользуясь неравенством $\max_{0 \leq i \leq N} |u_i| \leq \max_{0 \leq i \leq N} |y_i| + \max_{0 \leq i \leq N} |z_i|$ и оценками (15.27), (15.30), приходим к неравенству (15.31). ■

7. Устойчивость. Рассмотрим вопрос о чувствительности решения разностной схемы к погрешностям задания правых частей разностных уравнений. Пусть u^h — решение разностной схемы (15.15), (15.16), а u^{*h} — решение разностной схемы

$$L^h[u^{*h}] = f^{*h}, \quad x \in \omega^h; \quad (15.32)$$

$$u^{*h}(x_0) = u_a^*, \quad u^{*h}(x_N) = u_b^*, \quad (15.33)$$

где $f^{*h} = f^h - \delta f^h$, $u_a^* = u_a - \varepsilon_a$, $u_b^* = u_b - \varepsilon_b$.

Назовем разностную схему (15.15), (15.16) *устойчивой*, если при любых ε_a , ε_b , δf^h справедлива оценка

$$\max_{0 \leq i \leq N} |u_i - u_i^*| \leq \max\{|\varepsilon_a|, |\varepsilon_b|\} + K \max_{0 < i < N} |\delta f_i|, \quad (15.34)$$

где постоянная K не зависит от h . Отметим, что эта оценка является аналогом оценки (15.8), справедливой для краевой задачи (15.9), (15.10).

Теорема 15.13 (об устойчивости разностной схемы). Для разностной схемы (15.15), (15.16) справедлива оценка (15.34) с постоянной $K = l^2/8$.

Доказательство. Заметим, что сеточная функция $\varepsilon^h = u^h - u^{*h}$ является решением разностной схемы

$$L^h[\varepsilon^h] = \delta f^h, \quad x \in \omega^h;$$

$$\varepsilon^h(x_0) = \varepsilon_a, \quad \varepsilon^h(x_N) = \varepsilon_b.$$

Применяя для оценивания ε^h теорему 15.12, приходим к неравенству (15.34). ■

Замечание. Значение $l^2/8$ постоянной K в неравенстве (15.34) для разностной схемы (15.15), (15.16) совпадает в силу замечания 3 с. 558 со значением соответствующей постоянной в неравенстве (15.8). Этот факт говорит о том, что разностная схема обладает такой же чувствительностью к погрешностям задания исходных данных, что и краевая задача.

8. Аппроксимация. Пусть $u(x)$ — решение дифференциального уравнения $L[u] = f$. Назовем сеточную функцию $\Psi^h = L^h[u] - f^h$ *погрешностью аппроксимации разностного уравнения*

$$L^h[u^h] = f^h, \quad x \in \omega^h. \quad (15.35)$$

Из определения ψ^h следует, что справедливо равенство

$$L^h[u] = f^h + \psi^h, \quad x \in \omega^h. \quad (15.36)$$

означающее, что функция u удовлетворяет разностному уравнению (15.35) с точностью до погрешности аппроксимации.

Сеточную функцию ψ^h используют для предварительной оценки того, насколько точно аппроксимируется дифференциальное уравнение его разностным аналогом. Говорят, что *разностное уравнение* (15.35) *аппроксирует дифференциальное уравнение* $L[u] = f$, если $\max_{0 < i < N} |\psi_i| \rightarrow 0$ при $h \rightarrow 0$, и *аппроксирует его с m -м порядком* (при $m > 0$), если справедлива оценка $\max_{0 < i < N} |\psi_i| \leq Ch^m$.

Лемма 15.4. Пусть коэффициенты q и f дважды непрерывно дифференцируемы на отрезке $[a, b]$. Тогда разностное уравнение (15.15) аппроксимирует дифференциальное уравнение $L[u] = f$ со вторым порядком, причем справедлива оценка

$$\max_{0 < i < N} |\psi_i| \leq \frac{M_4}{12} h^2, \quad M_4 = \max_{[a, b]} |u^{(4)}(x)|. \quad (15.37)$$

Доказательство. Прежде всего заметим, что в силу теоремы 15.2 функция $u(x)$ имеет на отрезке $[a, b]$ непрерывную производную $u^{(4)}(x)$.

В силу определения погрешности аппроксимации имеем

$$\begin{aligned} \psi_i &= -\frac{u(x_{i-1}) - 2u(x_i) + u(x_{i+1})}{h^2} + q(x_i)u(x_i) - f(x_i) = \\ &= -u''(x_i) + r_h(x_i) + q(x_i)u(x_i) - f(x_i) = r_h(x_i), \end{aligned}$$

где $r_h(x_i) = u^{(4)}(\xi_i)h^2/12$ — погрешность аппроксимации производной $u''(x)$ ее разностным аналогом по формуле (15.11).

Таким образом, $\max_{0 < i < N} |\psi_i| = \max_{0 < i < N} \frac{|u^{(4)}(\xi_i)|}{12} h^2 \leq \frac{M_4}{12} h^2$. ■

В общем случае дополнительно возникает проблема изучения погрешности аппроксимации краевых условий. Однако для краевых условий (15.10) эта проблема отсутствует, так как сеточная функция удовлетворяет им точно.

9. Сходимость. Пусть $u(x)$ — решение краевой задачи, а u^h — решение соответствующей разностной схемы. *Назовем погрешностью*

разностной схемы сеточную функцию ε^h , принимающую значения $\varepsilon_i = \varepsilon^h(x_i) = u(x_i) - u^h(x_i)$ в узлах сетки $\bar{\omega}^h$.

Будем говорить, что *разностная схема сходится* при $h \rightarrow 0$, если $\max_{0 \leq i \leq N} |\varepsilon_i| \rightarrow 0$ при $h \rightarrow 0$, и *сходится с m -м порядком точности* (при $m > 0$), если для погрешности справедлива оценка $\max_{0 \leq i \leq N} |\varepsilon_i| \leq Ch^m$, где C — некоторая постоянная, не зависящая от h .

Покажем, что разностная схема (15.15), (15.16) сходится со вторым порядком точности.

Теорема 15.14. Пусть коэффициенты q и f дважды непрерывно дифференцируемы на отрезке $[a, b]$. Тогда для погрешности разностной схемы (15.15), (15.16) справедлива оценка

$$\max_{0 \leq i \leq N} |u(x_i) - u_i| \leq Ch^2, \quad (15.38)$$

где $C = \frac{l^2}{96} \max_{[a, b]} |u^{(4)}(x)|$.

Доказательство. Введем сеточную функцию u^{*h} , значения которой в узлах сетки совпадают с точными значениями решения краевой задачи, т.е. $u^{*h}(x_i) = u(x_i)$. Равенство (15.36) означает, что u^{*h} можно рассматривать как решение разностной схемы (15.32), (15.33), где $f^{*h} = f^h + \psi^h$, $u_a^* = u_a$, $u_b^* = u_b$. В силу теоремы 15.13 для $\varepsilon^h = u^h - u^{*h}$ справедлива оценка

$$\max_{0 \leq i \leq N} |\varepsilon_i| \leq \frac{l^2}{8} \max_{0 < i < N} |\psi_i|. \quad (15.39)$$

Учитывая теперь неравенство (15.37), из (15.39) получаем оценку (15.38). ■

Замечание. Сходимость разностной схемы (15.15), (15.16) со вторым порядком точности вытекает из того, что схема устойчива и обладает аппроксимацией со вторым порядком относительно h .

10. Оценка погрешности до правилу Рунге. Полученная в теореме 15.14 априорная оценка (15.38), как правило, оказывается непригодной для практической оценки погрешности разностной схемы. На практике чаще применяются апостериорные оценки погрешности, использующие расчеты на сгущающихся сетках. Пусть, например u^h и u^{2h} — решения разностной схемы (15.15), (15.16), соответствующие

шагам $h_1 = h$ и $h_2 = 2h$. Тогда в соответствии с правилом Рунге при определенных условиях справедлива приближенная формула

$$\varepsilon^h(x) = u(x) - u^h(x) \approx \frac{u^h(x) - u^{2h}(x)}{3}, \quad x \in \omega^{2h}. \quad (15.40)$$

Отметим, что она применима только в узлах сетки ω^{2h} , т.е. там, где определены обе сеточные функции u^h и u^{2h} .

Пример 15.1. Используя разностную схему (15.15), (15.16) с шагом $h = 1/8$, найдем приближенное решение краевой задачи

$$-u''(x) + x^2 u(x) = \left(\frac{\pi^2}{4} + x^2\right) \cos \frac{\pi}{2} x, \quad 0 < x < 1; \quad (15.41)$$

$$u(0) = 1, \quad u(1) = 0$$

и оценим его погрешность по правилу Рунге. Вычисления будем вести с шестью значащими цифрами.

В данном случае $q(x) = x^2$, $f(x) = \left(\frac{\pi^2}{4} + x^2\right) \cos \frac{\pi}{2} x$ и система сеточных уравнений (15.17), (15.18) принимает вид

$$u_0 = 1,$$

$$-u_{i-1} + (2 + h^2 x_i^2) u_i - u_{i+1} = \left(\frac{\pi^2}{4} + x_i^2\right) \cos \frac{\pi}{2} x_i, \quad 1 \leq i \leq N-1; \quad (15.42)$$

$$u_N = 0.$$

При $h_1 = 1/8$ приходим к следующей системе уравнений относительно неизвестных $u_i \approx u(ih)$ ($i = 0, 1, \dots, 8$):

$$u_0 = 1,$$

$$-u_0 + 2.00024u_1 - u_2 = 0.0380518,$$

$$-u_1 + 2.00098u_2 - u_3 = 0.0365207,$$

$$-u_2 + 2.00220u_3 - u_4 = 0.0338827,$$

$$-u_3 + 2.00391u_4 - u_5 = 0.0300233,$$

$$-u_4 + 2.00610u_5 - u_6 = 0.0248099,$$

$$-u_5 + 2.00897u_6 - u_7 = 0.0181171,$$

$$-u_6 + 2.01196u_7 - u_8 = 0.0098552,$$

$$u_6 = 0.$$

Таблица 15.1

x_i	$u = \cos \frac{\pi}{2} x$	$u^{h_1}, h_1 = 1/8$	$\varepsilon^{h_1}, h_1 = 1/8$	$u^{h_2}, h_2 = 1/4$	Оценка погрешности по правилу Рунге
0.000	1.000000	1.000000	0	1.000000	0
0.125	0.980785	0.981114	$-3 \cdot 10^{-4}$	—	—
0.250	0.923880	0.924413	$-5 \cdot 10^{-4}$	0.926080	$-6 \cdot 10^{-4}$
0.375	0.831470	0.832097	$-6 \cdot 10^{-4}$	—	—
0.500	0.707107	0.707733	$-6 \cdot 10^{-4}$	0.709703	$-7 \cdot 10^{-4}$
0.625	0.555570	0.556116	$-5 \cdot 10^{-4}$	—	—
0.750	0.382683	0.383082	$-4 \cdot 10^{-4}$	0.384324	$-4 \cdot 10^{-4}$
0.875	0.195090	0.195300	$-2 \cdot 10^{-4}$	—	—
1.000	0.000000	0.000000	0	0.000000	0

Решая ее с помощью метода прогонки, находим значения, представленные в третьем столбце табл. 15.1. Заметим, что в данном случае точное решение задачи известно. $u(x) = \cos \frac{\pi}{2} x$. Значения точного решения и вычисленные с их использованием погрешности приведены во втором и четвертом столбцах той же табл. 15.1.

Оценим теперь погрешность, использовав правило Рунге. Возьмем шаг $h_2 = 2h_1 = 1/4$. Решая соответствующую систему сеточных уравнений:

$$u_0 = 1,$$

$$-u_0 + 2.00391u_1 - u_2 = 0.146083,$$

$$-u_1 + 2.01563u_2 - u_3 = 0.120093,$$

$$-u_2 + 2.03516u_3 - u_4 = 0.0724683,$$

$$u_4 = 0,$$

получаем значения, представленные в пятом столбце табл. 15.1. В последнем столбце этой таблицы даны приближенные значения погрешности, полученные по формуле (15.40). Заметим, что в данном примере правило Рунге дает весьма хорошие результаты. ▲

11. Влияние вычислительной погрешности. При расчетах с достаточно крупными шагами h влиянием вычислительной погрешности на решение часто можно пренебречь. Однако все же следует иметь в виду, что при решении системы (15.17), (15.18) методом прогонки

происходит накопление вычислительной погрешности. Известно¹, что при $h \rightarrow 0$ вычислительная погрешность может возрастать здесь пропорционально ϵ_m/h^2 , где ϵ_m — относительная точность представления чисел в компьютере. Таким образом, при достаточно малых значениях шага h возможна катастрофическая потеря точности.

Сделанное выше утверждение перестанет казаться неправдоподобным, если мы убедимся в том, что при малых h вычислительная погрешность может привести к существенному искажению решения уже на этапе составления системы сеточных уравнений. Пусть, например для решения краевой задачи (15.9), (15.10) используется разностная схема с шагом $h = 10^{-3}$, а вычисления ведутся на 6-разрядном десятичном компьютере. Напомним, что для такого компьютера $\epsilon_m = 5 \cdot 10^{-7}$.

Так как здесь $h^2 x_i^2 < 10^{-6} = 2\epsilon_m$, то результатом вычисления коэффициента $b_i = 2 + h^2 x_i^2$ после округления до шести значащих цифр мантиссы является число $b_i^* = 2$. Следовательно, даже если остальные вычисления будут производиться точно, то фактически окажется найденным не решение системы (15.42), а решение системы

$$u_0^* = 1,$$

$$-u_{i-1}^* + 2u_i^* - u_{i+1}^* = h^2 \left(\frac{\pi^2}{4} + x_i^2 \right) \cos \frac{\pi}{2} x_i, \quad 1 \leq i \leq N-1;$$

$$u_N^* = 0,$$

соответствующей краевой задаче

$$-u''(x) = \left(\frac{\pi^2}{4} + x^2 \right) \cos \frac{\pi}{2} x;$$

$$u(0) = 1, \quad u(1) = 0.$$

Поскольку в результате оказалась решенной «не та задача», найденные значения u_i^* будут существенно отличаться от искомых u_i . В данном конкретном случае погрешность будет достигать примерно 5 %. В общем случае погрешность может оказаться значительно больше.

¹ Бахвалов Н.С. О накоплении вычислительной погрешности при численном решении дифференциальных уравнений. Вычислительные методы и программирование. М.: Изд-во МГУ, 1962. С. 47—68.

§ 15.3. Метод конечных разностей: аппроксимации специального вида

1. Случай переменного коэффициента $k(x)$. Вернемся к проблеме численного решения краевой задачи

$$-(k(x)u'(x))' + q(x)u(x) = f(x), \quad a < x < b; \quad (15.43)$$

$$u(a) = u_a \quad u(b) = u_b. \quad (15.44)$$

По сравнению со случаем $k \equiv 1$, рассмотренным в § 15.2, единственное видимое отличие состоит в необходимости выбора подходящей аппроксимации для выражения¹ $-(k(x)u'(x))' = w'(x)$. Рассмотрим некоторые из возможных подходов к выбору аппроксимации.

Введем обозначения $x_{i+1/2} = (x_i + x_{i+1})/2$, $k_{i+1/2} = k(x_{i+1/2})$, $0 \leq i < N$. Аппроксимируем производную $w'(x)$ при $x = x_i$ следующим образом:

$$w'(x_i) \approx \frac{1}{h} (w(x_{i+1/2}) - w(x_{i-1/2})). \quad (15.45)$$

Используя далее приближенные формулы

$$w(x_{i+1/2}) = -k_{i+1/2}u'(x_{i+1/2}) \approx -k_{i+1/2} \frac{u(x_{i+1}) - u(x_i)}{h}, \quad (15.46)$$

$$w(x_{i-1/2}) = -k_{i-1/2}u'(x_{i-1/2}) \approx -k_{i-1/2} \frac{u(x_i) - u(x_{i-1})}{h}, \quad (15.47)$$

получаем аппроксимацию

$$-(ku')'|_{x=x_i} \approx -\frac{1}{h} \left[k_{i+1/2} \frac{u(x_{i+1}) - u(x_i)}{h} - k_{i-1/2} \frac{u(x_i) - u(x_{i-1})}{h} \right].$$

В результате приходим к разностной схеме вида (15.15), (15.16), где

$$L^h[u^h](x_i) = -\frac{1}{h} \left[k_{i+1/2} \frac{u_{i+1} - u_i}{h} - k_{i-1/2} \frac{u_i - u_{i-1}}{h} \right] + q_i u_i. \quad (15.48)$$

Отметим, что коэффициенты $a_i = -k_{i-1/2}$, $b_i = k_{i-1/2} + k_{i+1/2} + h^2 q_i$, $c_i = -k_{i+1/2}$ соответствующей системы сеточных уравнений удовлетворяют условиям (15.23). Следовательно, при любом h решение разностной схемы существует и единственno. Кроме того, в силу леммы 15.1 для разностной схемы справедлив принцип максимума. Можно показать также, что разностная схема устойчива и сходится со вторым

¹ Напомним, что величина $w(x) = -k(x)u'(x)$ имеет физический смысл потока тепла, если уравнение интерпретируется как уравнение теплопроводности.

порядком точности, если коэффициенты k, q, f являются дважды непрерывно дифференцируемыми на отрезке $[a, b]$ функциями.

Формальный подход к выбору аппроксимации дифференциального уравнения может давать разностные схемы, обладающие теми или иными дефектами. Например, кажется удобным предварительно преобразовать первое слагаемое уравнения (15.43) следующим образом: $(ku')' = ku'' + k'u'$. После такого преобразования для этого слагаемого естествен выбор аппроксимации

$$(ku')'|_{x=x_i} \approx k(x_i) \frac{u(x_{i-1}) - 2u(x_i) + u(x_{i+1})}{h^2} + k'(x_i) \frac{u(x_{i+1}) - u(x_{i-1})}{2h},$$

приводящей к разностной схеме с оператором:

$$L^h[u^h](x_i) = -k(x_i) \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} - k'(x_i) \frac{u_{i+1} - u_{i-1}}{2h} + q_i u_i.$$

Заметим, что коэффициенты соответствующей системы сеточных уравнений

$$a_i = -k(x_i) + \frac{h}{2} k'(x_i), \quad b_i = 2k(x_i) + h^2 q_i, \quad c_i = -k(x_i) - \frac{h}{2} k'(x_i)$$

удовлетворяют условиям (15.23), гарантирующим наличие принципа максимума, только если

$$h \cdot \max_{1 \leq i < N} \frac{|k'(x_i)|}{|k(x_i)|} < 2. \quad (15.49)$$

Таким образом, в ситуации, когда коэффициент k может резко меняться на отрезке $[a, b]$, ограничение (15.49) приводит к необходимости выбора очень мелкого шага h для получения приемлемых результатов.

2. Случай неравномерной сетки. Часто возникает необходимость использования *неравномерной сетки* $\bar{\omega}^h$, т.е. сетки, у которой шаг $h_i = x_i - x_{i-1}$ зависит от i . Положим $h_{i+1/2} = x_{i+1/2} - x_{i-1/2}$. Заменив в формулах (15.45), (15.46), (15.47) h на $h_{i+1/2}, h_{i+1}, h_i$ соответственно, придем к разностной схеме (15.15), (15.16), в которой

$$L^h[u^h](x_i) = -\frac{1}{h_{i+1/2}} \left[k_{i+1/2} \frac{u_{i+1} - u_i}{h_{i+1}} - k_{i-1/2} \frac{u_i - u_{i-1}}{h_i} \right] + q_i u_i.$$

Нетрудно убедиться в том, что при такой аппроксимации справедлив принцип максимума и разностная схема устойчива. Можно показать, что при некоторых дополнительных предположениях она сходится со вторым порядком точности относительно h_{\max} .

3. Разности «против потока». Как отмечалось выше, уравнение (15.43) описывает установившееся распределение температуры в неподвижной среде. Когда исследуются тепловые процессы в движущейся среде (например, рассматривается поток жидкости), уравнение модифицируется следующим образом:

$$-(k(x) u'(x))' + v(x) u'(x) + q(x) u(x) = f(x).$$

Здесь $v(x)$ — величина, пропорциональная скорости потока жидкости.

При дискретизации этого уравнения возникает новый момент, связанный с необходимостью аппроксимации слагаемого $v(x) u'(x)$. Кажется естественным воспользоваться для аппроксимации производной u' центральной разностной производной. В результате к разностному

оператору (15.48) добавится слагаемое $v_i \frac{u_{i+1} - u_{i-1}}{2h}$, где $v_i = v(x_i)$.

Выясним, удовлетворяют ли коэффициенты

$$a_i = -k_{i-1/2} - \frac{h}{2} v_i, \quad b_i = k_{i-1/2} + k_{i+1/2} + h^2 q_i, \quad c_i = -k_{i+1/2} + \frac{h}{2} v_i$$

соответствующей системы сеточных уравнений условиям (15.23), гарантирующим выполнение принципа максимума. Как нетрудно видеть, неравенства (15.23) выполняются, если $h|v_i| < 2 \min\{k_{i-1/2}, k_{i+1/2}\}$.

Когда скорость потока велика, это неравенство приводит к весьма жесткому ограничению на шаг h . Его можно избежать, если использовать односторонние разностные производные. В задачах динамики жидкостей и газов широко используются аппроксимации вида

$$v(x_i) u'(x_i) \approx v_i^+ \frac{u(x_i) - u(x_{i-1})}{h} + v_i^- \frac{u(x_{i+1}) - u(x_i)}{h}, \quad (15.50)$$

где $v_i^- = \min\{v(x_i), 0\}$, $v_i^+ = \max\{v(x_i), 0\}$.

Они называются *аппроксимациями «против потока»* (или *«против ветра»*). Такой выбор аппроксимации слагаемого $v u'$ приводит к системе сеточных уравнений с коэффициентами

$$a_i = -(k_{i-1/2} + h v_i^+), \quad b_i = k_{i-1/2} + k_{i+1/2} + h v_i + h^2 q_i, \quad c_i = -(k_{i+1/2} - h v_i^-).$$

Легко убедиться в том, что условия (15.23) здесь всегда выполняются. Таким образом, принцип максимума выполняется при любых шагах h . Правда, при использовании приближения (15.50) порядок аппроксимации снижается со второго до первого.

4. Случай разрывных коэффициентов. Одна из специфических особенностей, присущих многим техническим задачам, заключается в том, что среда, в которой изучаются те или иные процессы, как правило, существенно неоднородна и состоит из материалов с разными физическими свойствами. При математической формулировке таких

задач эта особенность проявляется в том, что коэффициенты дифференциальных уравнений становятся разрывными. Это существенно усложняет построение эффективных численных методов.

Предположим, например, что коэффициенты k, q, f , входящие в уравнение (15.43), могут иметь на отрезке $[a, b]$ конечное число M точек разрыва ξ_i ($a < \xi_1 < \xi_2 < \dots < \xi_M < b$) первого рода. Будем предполагать, что всюду за исключением этих точек коэффициенты k, q, f непрерывны и удовлетворяют условиям $k(x) \geq k_0 > 0, q(x) \geq 0$.

В этом случае решение краевой задачи (15.43), (15.44) уже нельзя понимать в классическом смысле. Уточним постановку задачи для уравнения с разрывными коэффициентами. Назовем функцию $u(x)$ *решением задачи* (15.43), (15.44), если:

1) функция $u(x)$ непрерывна на отрезке $[a, b]$ и удовлетворяет краевым условиям $u(a) = u_a, u(b) = u_b$;

2) поток $w(x) = -k(x)u'(x)$ непрерывен на отрезке $[a, b]$;

3) всюду за исключением точек $\xi_1, \xi_2, \dots, \xi_M$ функция $w(x)$ непрерывно дифференцируема и удовлетворяет уравнению

$$w'(x) + q(x)u(x) = f(x). \quad (15.51)$$

Для вывода разностных уравнений воспользуемся *методом баланса*¹. Запишем уравнение теплового баланса для отрезка $[x_{i-1/2}, x_{i+1/2}]$, где $1 \leq i \leq N-1$. Для этого проинтегрируем уравнение (15.51) по x от $x_{i-1/2}$ до $x_{i+1/2}$. В результате получим равенство

$$w(x_{i+1/2}) - w(x_{i-1/2}) + \int_{x_{i-1/2}}^{x_{i+1/2}} qu \, dx = \int_{x_{i-1/2}}^{x_{i+1/2}} f \, dx. \quad (15.52)$$

Воспользовавшись приближенной формулой $\int q u \, dx \approx u(x_i) \int q \, dx$ и разделив обе части равенства (15.52) на $h_{i+1/2}$, получим

$$\frac{1}{h_{i+1/2}} (w(x_{i+1/2}) - w(x_{i-1/2})) + q_i^h u(x_i) \approx f_i^h. \quad (15.53)$$

Здесь $q_i^h = \frac{1}{h_{i+1/2}} \int_{x_{i-1/2}}^{x_{i+1/2}} q(x) \, dx$, $f_i^h = \frac{1}{h_{i+1/2}} \int_{x_{i-1/2}}^{x_{i+1/2}} f(x) \, dx$ — средние значения функций q и f на отрезке $[x_{i-1/2}, x_{i+1/2}]$.

¹ Иногда метод баланса называют *интегро-интерполяционным методом* [85, 86].

Далее заметим, что

$$u(x_{i+1}) - u(x_i) = \int_{x_i}^{x_{i+1}} u'(x) dx = - \int_{x_i}^{x_{i+1}} \frac{w(x)}{k(x)} dx \approx -w(x_{i+1/2}) \int_{x_i}^{x_{i+1}} \frac{dx}{k(x)}.$$

Таким образом,

$$w(x_{i+1/2}) \approx -k_{i+1/2}^h \frac{u(x_{i+1}) - u(x_i)}{h_{i+1}}, \quad (15.54)$$

где $k_{i+1/2}^h = \left[\frac{1}{h_{i+1}} \int_{x_i}^{x_{i+1}} \frac{dx}{k(x)} \right]^{-1}$ — эффективное значение коэффициента

теплопроводности на отрезке $[x_i, x_{i+1}]$. Заметим (это важно!), что усредняется фактически не коэффициент теплопроводности $k(x)$, а обратный к нему коэффициент теплового сопротивления $(k(x))^{-1}$.

Перейдем теперь от приближенных равенств (15.53), (15.54) к разностному уравнению

$$-\frac{1}{h_{i+1/2}} \left[k_{i+1/2}^h \frac{u_{i+1} - u_i}{h_{i+1}} - k_{i-1/2}^h \frac{u_i - u_{i-1}}{h_i} \right] + q_i^h u_i = f_i^h, \quad 1 \leq i \leq N-1. \quad (15.55)$$

Добавляя к (15.55) уравнения

$$u_0 = u_a, \quad u_N = u_b, \quad (15.56)$$

приходим к разностной схеме (15.55), (15.56).

Замечание. Разностные уравнения (15.55) записываются еди-нообразно во всех внутренних узлах сетки независимо от того, где расположены точки разрыва коэффициентов дифференциального уравнения. Это означает, что рассматриваемая разностная схема относится к классу однородных разностных схем [86].

5. Аппроксимация краевых условий. Выше при аппроксимации краевой задачи краевые условия первого рода $u(a) = u_a$, $u(b) = u_b$ не вызывали каких-либо затруднений и потому основное внимание уделялось аппроксимации дифференциального оператора. Однако краевые условия могут иметь более сложный вид и тогда возникает проблема их аппроксимации.

Рассмотрим, например краевое условие второго рода

$$-k(a)u'(a) = w_a. \quad (15.57)$$

Простейший подход к его аппроксимации состоит в замене производной $u'(a)$ разностным отношением $\frac{u(a+h) - u(a)}{h}$. В результате получается разностное краевое условие

$$-k(a) \frac{u_1 - u_0}{h} = w_a. \quad (15.58)$$

Так как

$$\frac{u(a+h) - u(a)}{h} = u'(a) + \frac{u''(a)}{2} h + \frac{u^{(3)}(a)}{6} h^2 + \dots,$$

то разностное уравнение (15.58) аппроксимирует краевое условие (15.57) лишь с первым порядком относительно h , что приводит к понижению порядка точности разностной схемы. Порядок аппроксимации краевого условия можно повысить разными способами. Например, можно заметить, что в силу дифференциального уравнения (15.43) при $x = a$ для решения u справедливо равенство

$$k(a) u''(a) = -k'(a) u'(a) + q(a) u(a) - f(a).$$

Таким образом,

$$\begin{aligned} -k(a) \frac{u(a+h) - u(a)}{h} &= w_a + \frac{h}{2} \left[-\frac{k'(a)}{k(a)} w_a - q(a) u(a) + f(a) \right] \\ &\quad - k(a) \frac{u^{(3)}(a)}{6} h^2 + \dots, \end{aligned}$$

и мы приходим к разностному краевому условию

$$-k(a) \frac{u_1 - u_0}{h} + \frac{h}{2} q(a) u_0 = \left[1 - \frac{h}{2} \frac{k'(a)}{k(a)} \right] w_a + \frac{h}{2} f(a),$$

аппроксимирующему краевое условие (15.57) со вторым порядком.

Другая аппроксимация второго порядка относительно h получится, если использовать метод баланса. Проинтегрируем уравнение (15.51) по x от x_0 до $x_{1/2} = x_0 + h/2$. В результате, учитывая, что $w(x_0) = w_a$, полу-

чаем равенство $w(x_{1/2}) - w_a + \int\limits_{x_0}^{x_{1/2}} q u \, dx = \int\limits_{x_0}^{x_{1/2}} f \, dx$. Откуда, используя при-

ближенное равенство (15.54), приходим к разностному уравнению

$$-k_{1/2}^h \frac{u_1 - u_0}{h} + \frac{h}{2} q_0^h u_0 = w_a + \frac{h}{2} f_0^h,$$

где $q_0^h = \frac{2}{h} \int\limits_{x_0}^{x_{1/2}} q \, dx$, $f_0^h = \frac{2}{h} \int\limits_{x_0}^{x_{1/2}} f \, dx$.

§ 15.4. Понятие о проекционных и проекционно-разностных методах. Методы Ритца и Галеркина. Метод конечных элементов

Наряду с методом конечных разностей значительной популярностью пользуются проекционные методы Ритца и Галеркина, а точнее — их современные варианты, объединяемые названием «метод конечных элементов» или «проекционно-сеточные методы».

1. Вариационная постановка краевой задачи. Вариационные методы, представляющие собой частный случай проекционных методов, используются для решения самых разнообразных задач на протяжении многих десятков лет. Эти методы применяются для решения тех задач физики и техники, которые могут быть описаны с помощью так называемых *вариационных принципов*. В соответствии с одним из простейших вариационных принципов функция $u(x)$, являющаяся решением задачи, должна быть стационарной точкой *вариационного функционала*

$$J(u) = \int_a^b F(x, u, u') dx. \quad (15.59)$$

Вариационный функционал, как правило, имеет определенный физический смысл. Нередко он выражает потенциальную энергию физической системы.

Обозначим через U множество функций, на котором определен функционал $J(u)$. Будем считать, что входящие в U функции удовлетворяют условиям

$$u(a) = u_a, \quad u(b) = u_b, \quad (15.60)$$

где значения u_a и u_b фиксированы. Предположим также, что в множество U входят все непрерывные кусочно-гладкие функции, принимающие на концах отрезка $[a, b]$ значения (15.60).

Предположим, что непрерывно дифференцируемая функция $u(x)$ является стационарной точкой функционала (15.59). Тогда, как известно из курса вариационного исчисления [115], эта функция должна удовлетворять дифференциальному уравнению

$$-\frac{d}{dx} F'_{u'}(x, u, u') + F'_u(x, u, u') = 0, \quad (15.61)$$

которое принято называть *уравнением Эйлера* (или *уравнением Эйлера—Лагранжа*). Таким образом, решение вариационной задачи оказывается решением краевой задачи (15.61), (15.60). Более того, при некоторых условиях эти задачи оказываются эквивалентными и возни-

кает возможность решать определенный класс краевых задач, используя методы вариационного исчисления.

Рассмотрим теперь функционал

$$J(u) = \frac{1}{2} \int_a^b (k(u')^2 + qu^2) dx - \int_a^b fu dx, \quad (15.62)$$

где $k(x), q(x), f(x)$ — кусочно-непрерывные функции, удовлетворяющие условиям $k(x) \geq k_0 > 0$, $q(x) \geq 0$. Поставим вариационную задачу о поиске точки минимума функционала (15.62) на множестве U . Как нетрудно видеть, в рассматриваемом случае

$$F(x, u, u') = \frac{1}{2} (k(x)(u')^2 + q(x)u^2) - f(x)u,$$

и уравнение Эйлера принимает следующий вид:

$$-(ku')' + qu = f. \quad (15.63)$$

Можно доказать, что функция u является точкой минимума функционала (15.62), т.е. удовлетворяет условию

$$J(u) = \min_{v \in U} J(v), \quad (15.64)$$

тогда и только тогда, когда она является решением краевой задачи (15.63), (15.60).

Отметим одно достоинство вариационной постановки задачи (15.63), (15.60). Она исключает необходимость требования наличия у рассматриваемых функций второй производной и даже непрерывности первой производной. Это обстоятельство оказывается весьма ценным для многих приближенных методов.

2. Метод Ритца. Рассмотрим приближенный метод решения вариационной задачи о поиске точки минимума функционала $J(u)$ на множестве U . Будем искать приближенное решение u^N в виде следующей линейной комбинации:

$$u^N(x) = \sum_{j=0}^N \alpha_j \varphi_j(x). \quad (15.65)$$

Здесь $\varphi_0, \varphi_1, \dots, \varphi_N$ — некоторые фиксированные функции, которые далее мы будем называть *базисными*. Предполагается, что система базисных функций линейно независима и линейными комбинациями (15.65) при соответствующем выборе коэффициентов $\alpha_0, \alpha_1, \dots, \alpha_N$ можно аппроксимировать решение u с желаемой степенью точности.

Обозначим через U^N множество всех функций вида (15.65) (при фиксированных $\varphi_0, \varphi_1, \dots, \varphi_N$), удовлетворяющих условиям $u^N(a) = u_a$, $u^N(b) = u_b$. Предположим далее, что базисные функции удовлетворяют

следующим условиям: $\phi_0(a) = 1$, $\phi_j(a) = 0$ для всех $j \geq 1$; $\phi_N(b) = 1$, $\phi_j(b) = 0$ для всех $j < N$. Тогда, как нетрудно видеть, условия (15.60) для $u = u^N$ выполняются тогда и только тогда, когда справедливы равенства

$$\alpha_0 = u_a, \quad \alpha_N = u_b. \quad (15.66)$$

Согласно *методу Ритца*¹, приближенное решение u^N определяется как функция, минимизирующая функционал J на множестве U^N . Таким образом, по определению

$$J(u^N) = \min_{v \in U^N} J(v). \quad (15.67)$$

Заметим, что задача (15.67) представляет собой задачу минимизации функции многих переменных. В самом деле, величина $J(v) = J\left(\sum_{j=0}^N \alpha_j \phi_j\right)$

является функцией $N - 1$ переменных $\alpha_1, \alpha_2, \dots, \alpha_{N-1}$ (значения $\alpha_0 = u_a$, $\alpha_N = u_b$ фиксированы). Согласно необходимому условию экстремума, минимум этой функции достигается при тех значениях параметров $\alpha_1, \alpha_2, \dots, \alpha_{N-1}$, для которых выполняются равенства

$$\frac{\partial}{\partial \alpha_i} J\left(\sum_{j=0}^N \alpha_j \phi_j\right) = 0, \quad i = 1, 2, \dots, N-1. \quad (15.68)$$

Добавляя к этим равенствам условия (15.66), приходим к системе уравнений (15.68), (15.66), из которых можно определить значения коэффициентов α_j ($j = 0, 1, \dots, N$) и тем самым — приближение u^N .

Применим метод Ритца к решению краевой задачи для уравнения (15.63) с краевыми условиями первого рода. Для функционала (15.62) имеем

$$\begin{aligned} J\left(\sum_{j=0}^N \alpha_j \phi_j\right) &= \frac{1}{2} \int_a^b \left[k \left(\sum_{j=0}^N \alpha_j \phi'_j \right)^2 + q \left(\sum_{j=0}^N \alpha_j \phi_j \right)^2 \right] dx - \int_a^b f \sum_{j=0}^N \alpha_j \phi_j dx, \\ \frac{\partial}{\partial \alpha_i} J\left(\sum_{j=0}^N \alpha_j \phi_j\right) &= \int_a^b \left[k \left(\sum_{j=0}^N \alpha_j \phi'_j \right) \phi'_i + q \left(\sum_{j=0}^N \alpha_j \phi_j \right) \phi_i \right] dx - \int_a^b f \phi_i dx = \\ &= \sum_{j=0}^N \left[\int_a^b (k \phi'_i \phi'_j + q \phi_i \phi_j) dx \right] \alpha_j - \int_a^b f \phi_i dx. \end{aligned}$$

¹ Вальтер Ритц (1878—1909) — швейцарский физик и математик.

Система (15.68), (15.66) в данном случае превращается в систему линейных алгебраических уравнений

$$\sum_{j=0}^N a_{ij} \alpha_j = b_i, \quad i = 1, 2, \dots, N-1; \quad (15.69)$$

$$\alpha_0 = u_a, \quad \alpha_N = u_b, \quad (15.70)$$

$$\text{где } a_{ij} = \int_a^b (k \varphi'_i \varphi'_j + q \varphi_i \varphi_j) dx, \quad b_i = \int_a^b f \varphi_i dx.$$

Исключив переменные α_0 и α_N , систему (15.69), (15.70) можно свести к эквивалентной системе уравнений

$$A\alpha = d.$$

Здесь $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_{N-1})^T$, A — матрица порядка $N-1$ с элементами a_{ij} ($i, j = 1, 2, \dots, N-1$), $d = (d_1, d_2, \dots, d_{N-1})^T$, где $d_i = b_i - a_{i0}u_a - a_{iN}u_b$, $i = 1, 2, \dots, N-1$.

Отметим, что матрица A — симметричная и положительно определенная.

3. Проекционная постановка краевой задачи. Краевая задача

$$L[u] = f, \quad x \in [a, b]; \quad (15.71)$$

$$u(a) = u_a, \quad u(b) = u_b \quad (15.72)$$

допускает вариационную постановку тогда и только тогда, когда дифференциальное уравнение (15.71) является уравнением Эйлера для некоторого функционала J . Таким образом, этим свойством обладают далеко не все задачи. Например, краевая задача для уравнения

$$-(ku')' + vu' + qu = f \quad (15.73)$$

при $v(x) \neq 0$ не допускает классической вариационной постановки.

Приведем проекционную постановку краевой задачи, которая имеет место и в том случае, когда задача не может быть сформулирована как вариационная.

Будем называть *пробной функцией* всякую непрерывную на отрезке $[a, b]$ кусочно-дифференцируемую функцию $\varphi(x)$, обращающуюся в нуль при $x = a$, $x = b$. Множество всех пробных функций обозначим через Φ . Умножив уравнение (15.71) на произвольную функцию $\varphi \in \Phi$ и проинтегрировав полученное равенство по x от a до b , получим *интегральное тождество*¹

$$\int_a^b L[u] \varphi dx = \int_a^b f \varphi dx. \quad (15.74)$$

¹ Равенство (15.74) называют интегральным тождеством, подчеркивая тем самым, что оно выполняется для любой пробной функции φ .

Итак, если функция u является решением дифференциального уравнения (15.71), то она удовлетворяет интегральному тождеству (15.74). В то же время, как следует из основной леммы вариационного исчисления [115], если интегральное тождество (15.74) выполняется для любой пробной функции ϕ , то $L[u] = f$.

Таким образом, краевую задачу (15.71), (15.72) можно сформулировать в следующей проекционной постановке. Требуется найти такую функцию u , которая удовлетворяет интегральному тождеству (15.74) для произвольной пробной функции $\phi \in \Phi$ и для которой выполнены краевые условия (15.72).

Приведем в качестве примера интегральное тождество

$$\int_a^b (- (ku')' + vu' + qu) \phi \, dx = \int_a^b f \phi \, dx, \quad (15.75)$$

соответствующее дифференциальному уравнению (15.73). Придадим тождеству (15.75) несколько иную форму. Используя формулу интегрирования по частям, с учетом равенств $\phi(a) = 0$, $\phi(b) = 0$ получаем

$$-\int_a^b (ku')' \phi \, dx = -ku' \phi \Big|_a^b + \int_a^b ku' \phi' \, dx = \int_a^b ku' \phi' \, dx.$$

В результате интегральное тождество примет вид

$$\int_a^b (ku' \phi' + vu' \phi + qu\phi) \, dx = \int_a^b f \phi \, dx. \quad (15.76)$$

Отметим, что после замены уравнения (15.73) интегральным тождеством (15.76) возникает возможность рассматривать функции $u(x)$, обладающие лишь первыми производными. Это обстоятельство играет важную роль при построении и исследовании методов решения рассматриваемой задачи, а также ряда других задач. Кроме того, проекционная постановка оказывается удобной при рассмотрении уравнений с разрывными коэффициентами.

5. Метод Галеркина. Как и в методе Ритца, в методе Галеркина¹ приближенное решение ищется в виде

$$u^N(x) = \sum_{j=0}^N \alpha_j \varphi_j(x).$$

¹ Борис Григорьевич Галеркин (1871—1945) — русский инженер, ученый, специалист в области строительной механики и теории упругости.

Однако в отличие от метода Ритца основой для построения метода является не вариационная, а более общая проекционная постановка задачи. За приближенное решение в методе Галеркина принимается функция $u^N \in U^N$, которая удовлетворяет интегральному тождеству

$$\int_a^b L[u^N] \phi \, dx = \int_a^b f \phi \, dx$$

для любой пробной функции вида $\phi = \varphi_i$, $i = 1, 2, \dots, N - 1$.

Для задачи (15.71), (15.72) метод Галеркина с использованием интегрального тождества (15.76) приводит к следующей системе уравнений:

$$\begin{aligned} \int_a^b \left[k \left(\sum_{j=0}^N \alpha_j \varphi'_j \right) \varphi'_i + v \left(\sum_{j=0}^N \alpha_j \varphi'_j \right) \varphi_i + q \left(\sum_{j=0}^N \alpha_j \varphi_j \right) \varphi_i \right] dx &= \\ &= \int_a^b f \varphi_i \, dx, \quad i = 1, 2, \dots, N - 1; \\ \alpha_0 &= u_a, \quad \alpha_N = u_b. \end{aligned}$$

Заметим, что эта система при $v(x) \equiv 0$ в точности совпадает с системой (15.69), (15.70), полученной методом Ритца. Таким образом, применительно к решению краевой задачи (15.63), (15.60) методы Ритца и Галеркина оказываются эквивалентными.

Замечание. Приближенное решение u^h определяемое методом конечных разностей, задается только в узлах сетки $\bar{\omega}^h$, поэтому для получения значения решения в произвольной точке $x \in [a, b]$ приходится производить интерполяцию. В то же время проекционные методы дают в качестве приближенного решения функцию (15.65), вычисляемую в произвольной точке x .

Как мы отмечали в предыдущих параграфах, системы сеточных уравнений, получаемые методом конечных разностей, обладают тем важным свойством, что матрицы коэффициентов этих систем являются разреженными (более того, в рассмотренных примерах матрицы были трехдиагональными). Для решения таких систем разработаны эффективные методы.

В общем случае применение методов Ритца и Галеркина к решению краевых задач приводит к необходимости вычислять решения систем уравнений вида $A\alpha = d$ с заполненными (и часто плохо обусловленными) матрицами A . Современные варианты проекционных методов,

объединяемые термином «метод конечных элементов», свободны от указанного недостатка. Перейдем к их рассмотрению.

6. Метод конечных элементов. *Метод конечных элементов* представляет собой разновидность проекционных методов, основанную на специальном выборе базисных функций.

История метода весьма поучительна. Метод конечных элементов впервые был предложен Р. Курантом¹ в 1943 г., но тогда его важная работа опередила потребности практики и фактически осталась незамеченной. Затем в начале 50-х годов XX в. инженерами — специалистами по строительной механике был разработан новый подход к решению задач теории упругости. В тех случаях когда расчетная область имела сложную геометрию, она разбивалась на подобласти простой геометрии, в каждой из которых решение могло быть найдено аналитически. Эти подобласти были названы конечными элементами, а сам подход — методом конечных элементов. Только в начале 60-х годов математиками были осознаны практическое значение и математическая природа метода. В 60-х и 70-х годах шло бурное развитие теории метода, он завоевывал все более широкие области применения. К настоящему времени метод конечных элементов получил самое широкое распространение в вычислительной практике. На его основе разработано большое число пакетов прикладных программ для решения разнообразных инженерных и научных задач.

Отметим характерные черты метода конечных элементов, выделяющие его среди других проекционных методов.

1) Расчетная область (множество изменения независимой переменной) разбивается на конечное число элементарных подмножеств стандартной формы (которые и называют *конечными элементами*).

2) Используемые базисные функции φ_i таковы, что они:

на каждом элементе имеют простой вид (чаще всего — многочлены); отличны от нуля лишь на нескольких соседних элементах.

Покажем, как применяется метод конечных элементов к решению краевой задачи (15.63), (15.60). Разобьем отрезок $[a, b]$ точками $a = x_0 < x_1 < \dots < x_N = b$ на N элементарных отрезков $[x_{i-1}, x_i]$ длины h_i . Таким образом, здесь в роли конечного элемента выступает элементарный отрезок $[x_{i-1}, x_i]$.

¹ Рихард Курант (1888—1972) — немецкий математик.

Введем базисные функции $\varphi_j(x)$ для $j = 1, 2, \dots, N - 1$ следующим образом

$$\varphi_j(x) = \begin{cases} (x - x_{j-1})/h_j & \text{при } x \in [x_{j-1}, x_j], \\ (x_{j+1} - x)/h_{j+1} & \text{при } x \in [x_j, x_{j+1}], \\ 0 & \text{при } x \notin [x_{j-1}, x_{j+1}]. \end{cases} \quad (15.77)$$

График такой базисной функции («шапочки») изображен на рис. 15.2, а

Подчеркнем, что функция φ_j отлична от нуля только лишь на двух соседних конечных элементах (отрезках $[x_{j-1}, x_j]$ и $[x_j, x_{j+1}]$) и является кусочно-линейной.

Введем также функции $\varphi_j(x)$ для $j = 0$ и $j = N$ (рис. 15.2, б, в):

$$\varphi_0(x) = \begin{cases} (x_1 - x)/h_1 & \text{при } x \in [x_0, x_1], \\ 0 & \text{при } x \notin [x_0, x_1]; \end{cases} \quad (15.78)$$

$$\varphi_N(x) = \begin{cases} 0 & \text{при } x \notin [x_{N-1}, x_N], \\ (x - x_{N-1})/h_N & \text{при } x \in [x_{N-1}, x_N]. \end{cases} \quad (15.79)$$

Будем искать приближенное решение задачи в виде

$$u^N(x) = \sum_{j=0}^N \alpha_j \varphi_j(x) \quad (15.80)$$

Заметим, что базисные функции обладают тем свойством, что $\varphi_j(x_i) = 1$ и $\varphi_j(x_i) = 0$ при $i \neq j$. В силу этого $\alpha_j = u^N(x_j)$, т.е. коэффициенты линейной комбинации (15.80) представляют собой значения функ-

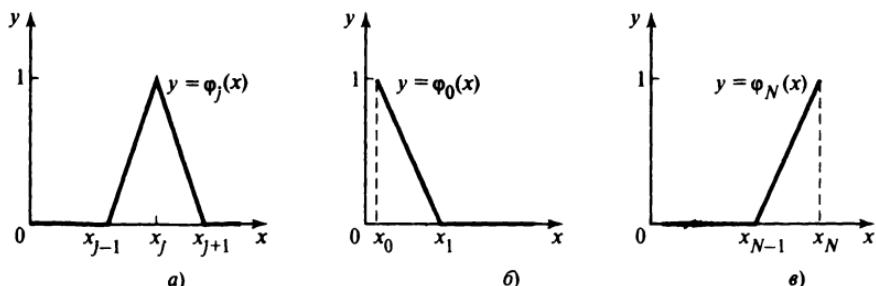


Рис. 15.2

ции u^N в узлах x_j . После введения обозначения $u_j^h = u^N(x_j)$ функцию (15.80) можно записать так

$$u^N(x) = \sum_{j=0}^N u_j^h \varphi_j(x).$$

Величины $u_j^h = a_j$ ($j = 0, 1, \dots, N$) удовлетворяют системе уравнений (15.69), (15.70), которую в данном случае можно получить как методом Ритца, так и методом Галеркина. Заметим, что базисные функции $\varphi_i(x)$ и $\varphi_j(x)$ одновременно могут быть отличны от нуля только, если $|i - j| \leq 1$,

поэтому при $|i - j| > 1$ элементы $a_{ij} = \int_a^b (k \varphi'_i \varphi'_j + q \varphi_i \varphi_j) dx$ равны

нулю. Таким образом, для определения неизвестных u_j^h получаем следующую систему уравнений с трехдиагональной матрицей.

$$a_{i,i-1}u_{i-1}^h + a_{i,i}u_i^h + a_{i,i+1}u_{i+1}^h = b_i, \quad 1 \leq i \leq N-1; \quad (15.81)$$

$$u_0^h = u_a, \quad u_N^h = u_b. \quad (15.82)$$

Так как

$$\varphi_j(x) = \begin{cases} 1/h_j & \text{при } x \in (x_{j-1}, x_j), \\ -1/h_{j+1} & \text{при } x \in (x_j, x_{j+1}), \\ 0 & \text{при } x \notin [x_{j-1}, x_{j+1}], \end{cases}$$

то

$$a_{i,i-1} = -h_{i-1}^{-1} k_{i-1/2}^h + h_{i+1/2} q_{i-1/2}^h, \quad (15.83)$$

$$a_{i,i+1} = -h_{i+1}^{-1} k_{i+1/2}^h + h_{i+1/2} q_{i+1/2}^h,$$

$$a_{i,i} = h_i^{-1} k_{i-1/2}^h + h_{i+1}^{-1} k_{i+1/2}^h + h_{i+1/2} q_i^h, \quad b_i = h_{i+1/2} f_i^h, \quad (15.84)$$

где

$$k_{i-1/2}^h = \frac{1}{h_i} \int_{x_{i-1}}^{x_i} k(x) dx, \quad q_{i\pm 1/2}^h = \frac{1}{h_{i\pm 1/2}} \int_{x_{i-1}}^{x_{i+1}} q(x) \varphi_i(x) \varphi_{i\pm 1}(x) dx, \quad (15.85)$$

$$q_i^h = \frac{1}{h_{i+1/2}} \int_{x_{i-1}}^{x_{i+1}} q(x) \varphi_i^2(x) dx, \quad f_i^h = \frac{1}{h_{i+1/2}} \int_{x_{i-1}}^{x_{i+1}} f(x) \varphi_i(x) dx. \quad (15.86)$$

Систему уравнений (15.81), (15.82) принято называть системой метода конечных элементов или *проекционно-разностной схемой*¹.

Можно показать, что проекционно-разностная схема имеет единственное решение $u^N(x)$ и при выполнении некоторых условий на сетку $\bar{\omega}^h$ и коэффициенты k, q, f она имеет второй порядок точности.

Существует весьма тесная связь между теорией разностных схем и теорией проекционно-разностных схем. В подтверждение сказанного ограничимся тем, что преобразуем систему (15.81), (15.82) так, чтобы она обрела внешнее сходство с соответствующей разностной схемой. Разделив каждое уравнение (15.81) на $h_{i+1/2}$ и воспользовавшись равенствами (15.83)–(15.87), получим систему уравнений

$$-\frac{1}{h_{i+1/2}} \left[k_{i+1/2}^h \frac{u_{i+1}^h - u_i^h}{h_{i+1}} - k_{i-1/2}^h \frac{u_i^h - u_{i-1}^h}{h_i} \right] +$$

$$+ q_{i-1/2}^h u_{i-1}^h + q_i^h u_i^h + q_{i+1/2}^h u_{i+1}^h = f_i^h \quad 1 \leq i \leq N-1; \quad (15.87)$$

$$u_0^h = u_a, \quad u_N^h = u_b. \quad (15.88)$$

В такой форме записи она действительно оказывается похожа на разностную схему (15.55), (15.56).

Замечание. При построении системы уравнений метода конечных элементов, как правило, возникает необходимость вычисления некоторых интегралов². Для проекционно-разностной схемы (15.87), (15.88) в случае, когда коэффициенты k, q, f — гладкие, эта проблема легко решается применением квадратурных формул. Например, можно положить

$$\int_{x_{i-1}}^{x_i} k(x) dx \approx h_i k(x_{i-1/2}), \quad \int_{x_{i-1}}^{x_i} q(x) \varphi_i(x) \varphi_{i-1}(x) dx \approx$$

$$\approx q(x_{i-1/2}) \int_{x_{i-1}}^{x_i} \frac{(x - x_{i-1})}{h_i} \frac{(x_i - x)}{h_i} dx = \frac{h_i}{6} q(x_{i-1/2}),$$

¹ Используется также термин *проекционно-сеточный метод*. Та же система уравнений может называться *вариационно-разностной схемой*, если она получена с помощью метода Ритца.

² Аналогичная проблема возникает и при реализации некоторых разностных схем, например однородной разностной схемы (15.55), (15.56).

$$\begin{aligned}
 & \int_{x_{i-1}}^{x_{i+1}} q(x) \varphi_i^2(x) dx \approx q(x_{i-1/2}) \int_{x_{i-1}}^{x_i} \frac{(x - x_{i-1})^2}{h_i^2} dx + \\
 & + q(x_{i+1/2}) \int_{x_i}^{x_{i+1}} \frac{(x_{i+1} - x)^2}{h_{i+1}^2} dx = \frac{h_i}{3} q(x_{i-1/2}) + \frac{h_{i+1}}{3} q(x_{i+1/2}); \\
 & \int_{x_i}^{x_{i+1}} q(x) \varphi_i(x) \varphi_{i+1}(x) dx \approx \frac{h_{i+1}}{6} q(x_{i+1/2}), \\
 & \int_{x_{i-1}}^{x_{i+1}} f(x) \varphi_i(x) dx \approx f(x_{i-1/2}) \int_{x_{i-1}}^{x_i} \frac{x - x_{i-1}}{h_i} dx + \\
 & + f(x_{i+1/2}) \int_{x_i}^{x_{i+1}} \frac{x_{i+1} - x}{h_{i+1}} dx = \frac{h_i}{2} f(x_{i-1/2}) + \frac{h_{i+1}}{2} f(x_{i+1/2}).
 \end{aligned}$$

При этом второй порядок точности сохраняется.

7. Специальная проекционно-разностная схема. Предположим, что коэффициенты k, q, f , входящие в одномерное уравнение диффузии (15.63), могут быть разрывными. В этом случае можно, как и ранее, использовать кусочно-линейные базисные функции (15.77), (15.78), (15.79). Полученная таким образом проекционно-разностная схема сходится, однако по скорости сходимости она существенно уступает однородной разностной схеме (15.55), (15.56).

Чтобы получить качественную проекционно-разностную схему для уравнения диффузии с разрывными коэффициентами, воспользуемся специальными базисными функциями $\varphi_j(x)$. Эти функции для $j = 1, 2, \dots, N - 1$ имеют следующий вид:

$$\varphi_j(x) = \begin{cases} \frac{1}{\Delta_j} \int_{x_{j-1}}^x \frac{dy}{k(y)} & \text{при } x \in [x_{j-1}, x_j], \\ \frac{1}{\Delta_{j+1}} \int_x^{x_{j+1}} \frac{dy}{k(y)} & \text{при } x \in [x_j, x_{j+1}], \\ 0 & \text{при } x \notin [x_{j-1}, x_{j+1}]. \end{cases} \quad (15.89)$$

Здесь $\Delta_j = \int_{x_{j-1}}^{x_j} \frac{dy}{k(y)}$, $\Delta_{j+1} = \int_{x_j}^{x_{j+1}} \frac{dy}{k(y)}$. Кроме того,

$$\varphi_0(x) = \begin{cases} \frac{1}{\Delta_1} \int_x^{x_1} \frac{dy}{k(y)} & \text{при } x \in [x_0, x_1], \\ 0 & \text{при } x \notin [x_0, x_1]; \end{cases} \quad (15.90)$$

$$\varphi_N(x) = \begin{cases} 0 & \text{при } x \notin [x_{N-1}, x_N], \\ \frac{1}{\Delta_N} \int_{x_{N-1}}^x \frac{dy}{k(y)} & \text{при } x \in [x_{N-1}, x_N]. \end{cases} \quad (15.91)$$

Выбранные базисные функции φ_i интересны тем, что всюду за исключением узлов сетки они удовлетворяют дифференциальному уравнению $(k\varphi'_i)' = 0$.

Применение этих базисных функций приводит к проекционно-разностной схеме вида (15.87), (15.88), где все коэффициенты находятся по формулам (15.85), (15.86) за исключением коэффициентов $k_{i-1/2}^h$, $k_{i+1/2}^h$, которые вычисляются по формулам

$$k_{i-1/2}^h = \left[\frac{1}{h_i} \int_{x_{i-1}}^{x_i} \frac{dx}{k(x)} \right]^{-1}, \quad k_{i+1/2}^h = \left[\frac{1}{h_{i+1}} \int_{x_i}^{x_{i+1}} \frac{dx}{k(x)} \right]^{-1}$$

т.е. так же, как и в разностной схеме (15.55), (15.56).

Эта проекционно-разностная схема имеет второй порядок точности и решается методом прогонки. Отметим одно замечательное свойство указанной схемы. Если $q(x) \equiv 0$, то полученные с ее помощью значения u_i^h совпадают с истинными значениями решения краевой задачи $u(x_i)$ в узлах сетки (схема точна в узлах сетки). Отметим, что базисные функции (15.89), (15.90), (15.91) и в случае, когда коэффициенты уравнения диффузии являются гладкими, но сильно меняющимися, дают большую точность по сравнению с простейшими кусочно-линейными базисными функциями.

§ 15.5. Метод пристрелки

Метод пристрелки (он называется также *методом стрельбы* или *баллистическим методом*) позволяет свести решение краевой задачи к решению системы нелинейных уравнений относительно так называемых пристрелочных параметров, а также к решению (вообще говоря, многократному) задачи Коши.

Сначала рассмотрим этот метод на примере решения следующей краевой задачи для системы двух дифференциальных уравнений первого порядка:

$$y'(x) = f(x, y(x), z(x)), \quad (15.92)$$

$$z'(x) = g(x, y(x), z(x)), \quad (15.93)$$

$$y(a) = y_a, \quad z(b) = z_b. \quad (15.94)$$

Наряду с этой задачей рассмотрим задачу Коши

$$y'(x, \alpha) = f(x, y(x, \alpha), z(x, \alpha)), \quad (15.95)$$

$$z'(x, \alpha) = g(x, y(x, \alpha), z(x, \alpha)), \quad (15.96)$$

$$y(a, \alpha) = y_a, \quad z(a, \alpha) = \alpha. \quad (15.97)$$

Решение задачи (15.95)–(15.97), т.е. пара функций $y(x, \alpha), z(x, \alpha)$, зависит не только от переменной x , но и от *пристрелочного параметра* α . Подобрав значение параметра α , при котором $z(b, \alpha) = z_b$, получим решение задачи Коши, совпадающее с решением краевой задачи (15.92)–(15.94).

Таким образом, для того чтобы найти решение краевой задачи (15.92)–(15.94) методом пристрелки, нужно решить нелинейное уравнение

$$\psi(\alpha) = 0, \quad (15.98)$$

где $\psi(\alpha) = z(b, \alpha) - z_b$. Отметим, что функция $\psi(\alpha)$ не задана какой-либо явной формулой и вычисление каждого ее значения предполагает вычисление решения задачи Коши (15.95)–(15.97). Как правило, для решения задачи Коши приходится использовать тот или иной численный метод.

Уравнение (15.98) можно решать, используя один из известных методов решения нелинейных уравнений. Довольно часто успешным оказывается применение метода бисекции, метода секущих или метода Ньютона (см. гл. 4).

Применение метода Ньютона

$$\alpha_{k+1} = \alpha_k - \psi(\alpha_k)/\psi'(\alpha_k) \quad (15.99)$$

сопряжено с необходимостью вычисления значений не только функции $\psi(\alpha)$, но и ее производной $\psi'(\alpha)$. Покажем, как можно вычислить $\psi'(\alpha) = \frac{\partial}{\partial \alpha} z(b, \alpha)$ в рассматриваемом случае.

Дифференцируя по параметру α уравнения (15.95), (15.96) и начальные условия (15.97), замечаем, что функции $u(x, \alpha) = \frac{\partial}{\partial \alpha} y(x, \alpha)$, $v(x, \alpha) = \frac{\partial}{\partial \alpha} z(x, \alpha)$, удовлетворяют следующим уравнениям и начальным условиям:

$$u' = f'_y(x, y(x, \alpha), z(x, \alpha))u + f'_z(x, y(x, \alpha), z(x, \alpha))v, \quad (15.100)$$

$$v' = g'_y(x, y(x, \alpha), z(x, \alpha))u + g'_z(x, y(x, \alpha), z(x, \alpha))v, \quad (15.101)$$

$$u(a, \alpha) = 0, \quad v(a, \alpha) = 1. \quad (15.102)$$

Решив теперь относительно функций y, z, u, v задачу Коши для системы уравнений (15.95), (15.96), (15.100), (15.101) с начальными условиями (15.97), (15.102), можно определить $\psi(\alpha) = z(b, \alpha) - z_b$ и $\psi'(\alpha) = v(b, \alpha)$.

Пример 15.2. Применим метод пристрелки к решению краевой задачи

$$y'(x) = \frac{x}{z(x)}, \quad z'(x) = -\frac{x}{y(x)}, \quad 0 \leq x \leq 1, \quad (15.103)$$

$$y(0) = 1, \quad z(1) = 1. \quad (15.104)$$

Положим $\psi(\alpha) = z(1, \alpha) - 1$, где $y(x, \alpha), z(x, \alpha)$ — решение задачи Коши

$$y'(x, \alpha) = \frac{x}{z(x, \alpha)}, \quad z'(x, \alpha) = -\frac{x}{y(x, \alpha)}, \quad (15.105)$$

$$y(0, \alpha) = 1, \quad z(0, \alpha) = \alpha. \quad (15.106)$$

Для решения уравнения $\psi(\alpha) = 0$ воспользуемся методом секущих:

$$\alpha_{k+1} = \alpha_k - \psi(\alpha_k) \frac{\alpha_k - \alpha_{k-1}}{\psi(\alpha_k) - \psi(\alpha_{k-1})}. \quad (15.107)$$

Возьмем $\alpha_0 = 0.5$, $\alpha_1 = 1$ и будем вести итерации по формуле (15.107) до тех пор, пока не выполнится условие $|\psi(\alpha_k)| < 10^{-6}$. Результаты вычислений приведены в табл. 15.2.

Таблица 15.2

k	α_k	$\psi(\alpha_k)$
0	0.500000	-0.816060
1	1.000000	-0.393469
2	1.46554	0.419069 · 10 ⁻¹
3	1.42073	-0.760640 · 10 ⁻³
4	1.42153	0.6098 · 10 ⁻⁷

Заметим, что в общем случае вычисление значения $\varphi(\alpha)$ производится с помощью численного решения задачи Коши. Однако в рассматриваемом случае задача (15.105), (15.106) допускает аналитическое решение

$$y(x, \alpha) = e^{\frac{x^2}{2\alpha}}, \quad z(x, \alpha) = \alpha e^{-\frac{x^2}{2\alpha}}. \quad (15.108)$$

Поэтому можно воспользоваться явной формулой

$$\psi(\alpha) = \alpha e^{-\frac{1}{2\alpha}} - 1.$$

В результате применения метода пристрелки для задачи (15.103), (15.104) получаем решение (15.108), где $\alpha \approx 1.42153$. ▲

Замечание. Своим названием метод пристрелки обязан очевидной аналогии между процессом его реализации при решении краевой задачи и процессом пристрелки при артиллерийской стрельбе по цели. После выбора очередного значения α_k пристрелочного параметра (выбора угла стрельбы) решается задача Коши (производится «выстрел»). Если $z(b, \alpha_k)$ совпадает с z_b заданной точностью, то «цель считается пораженной», а краевая задача — решенной. В противном случае производится корректировка значения пристрелочного параметра и процесс продолжается дальше. Использование для решения уравнения (15.98) метода бисекции еще более усиливает эту аналогию. Здесь результат того «выстрела», при котором $\psi(\alpha_k) > 0$, может восприниматься как «перелет снаряда», а того, при котором $\psi(\alpha_k) < 0$, — как «недолет».

Покажем теперь схематично, как применяется метод пристрелки для решения общей двухточечной краевой задачи для системы дифференциальных уравнений первого порядка

$$y'_1(x) = f_1(x, y_1(x), \dots, y_n(x)),$$

.....

$$y'_n(x) = f_n(x, y_1(x), \dots, y_n(x)),$$

$$\varphi_1(y_1(a), \dots, y_n(a), y_1(b), \dots, y_n(b)) = 0,$$

$$\varphi_n(y_1(a), \dots, y_n(a), y_1(b), \dots, y_n(b)) = 0.$$

Запишем эту задачу в векторной форме:

$$y'(x) = f(x, y(x)), \quad (15.109)$$

$$\varphi(y(a), y(b)) = 0. \quad (15.110)$$

Рассмотрим также задачу Коши

$$y'(x, \alpha) = f(x, y(x, \alpha)), \quad (15.111)$$

$$y(a, \alpha) = \alpha. \quad (15.112)$$

Здесь $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$ — вектор пристрелочных параметров.

Решив задачу Коши (15.111), (15.112) при фиксированном значении вектора α , получим решение $y(x, \alpha)$. Далее взяв $y(a) = \alpha$ и $y(b) = y(b, \alpha)$ и подставив эти значения в систему (15.110), придем к системе уравнений

$$\psi(\alpha) = 0, \quad (15.113)$$

где $\psi(\alpha) = \phi(\alpha, y(b, \alpha))$.

Наконец, решив систему (15.113), получим набор $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$ значений пристрелочных параметров, при которых решение задачи Коши (15.111), (15.112) совпадает с решением краевой задачи (15.109), (15.110).

Практическая реализация метода пристрелки при большом числе уравнений (часто уже при $n \geq 3$) оказывается довольно сложным делом. Действительно, даже сама по себе проблема решения системы нелинейных уравнений (15.113) является весьма трудной. Серьезные затруднения могут возникнуть здесь уже на этапе выбора хорошего начального приближения α_0 . Необходимо также учесть, что каждое вычисление вектор-функции $\psi(\alpha)$ является здесь весьма трудоемкой операцией: оно предполагает (численное) решение задачи Коши (15.111), (15.112).

Метод пристрелки достаточно эффективен в том случае, когда задача Коши (15.111), (15.112) является хорошо обусловленной. Однако если задача Коши плохо обусловлена, метод оказывается практически непригодным. Дело в том, что при решении системы (15.113) значения пристрелочных параметров α обязательно будут найдены с некоторой погрешностью, (относительное значение которой не может иметь порядок, меньший чем машинное эпсилон ϵ_m). Соответствующее решение задачи Коши (в случае плохой обусловленности) в результате этой погрешности окажется полностью искаженным. Однако даже в том идеализированном случае, когда вектор α найден абсолютно точно, при численном решении задачи (15.111), (15.112) на компьютере в приближенное решение будут внесены погрешности, которые сделают его непригодным. Для некоторых систем эти погрешности могут приводить даже к аварийному останову вычислительного процесса.

Пример 15.3. Рассмотрим плохо обусловленную краевую задачу

$$y'(x) = z(x), \quad z'(x) = 100y(x) + e^x, \quad 0 \leq x \leq 2; \quad (15.114)$$

$$y(0) = 0, \quad z(2) = 0. \quad (15.115)$$

Как нетрудно проверить, ее решением является пара функций

$$y(x) = c_1 e^{10x} + c_2 e^{-10x} - \frac{e^x}{99},$$

$$z(x) = 10c_1 e^{10x} - 10c_2 e^{-10x} - \frac{e^x}{99},$$

где

$$c_1 = \frac{1}{990} \frac{10e^{-40} + e^{-18}}{1 + e^{-40}} \approx 1.53838 \cdot 10^{-11},$$

$$c_2 = \frac{1}{990} \frac{10 - e^{-18}}{1 + e^{-40}} \approx 0.0101010.$$

Попробуем решить задачу (15.114), (15.115) методом пристрелки, использовав 6-разрядный десятичный компьютер. Соответствующая задача Коши имеет вид:

$$\begin{aligned} y'(x, \alpha) &= z(x, \alpha), \quad z'(x, \alpha) = 100y(x, \alpha) + e^x, \\ y(0, \alpha) &= 0, \quad z(0, \alpha) = \alpha. \end{aligned}$$

Ее решением являются функции

$$y(x, \alpha) = \frac{1}{20} \left(\frac{1}{9} + \alpha \right) e^{10x} + \frac{1}{20} \left(\frac{1}{11} - \alpha \right) e^{-10x} - \frac{1}{99} e^x, \quad (15.116)$$

$$z(x, \alpha) = \frac{1}{2} \left(\frac{1}{9} + \alpha \right) e^{10x} - \frac{1}{2} \left(\frac{1}{11} - \alpha \right) e^{-10x} - \frac{1}{99} e^x. \quad (15.117)$$

Уравнение $\psi(\alpha) = 0$ (где $\psi(\alpha) \equiv z(2, \alpha)$) для определения пристрелочного параметра α является линейным, поэтому для определения α достаточно сделать одну итерацию метода секущих:

$$\alpha = \alpha_1 - \frac{\psi(\alpha_1)(\alpha_1 - \alpha_0)}{\psi(\alpha_1) - \psi(\alpha_0)}. \quad (15.118)$$

Возьмем $\alpha_0 = 0$, $\alpha_1 = -1$. Тогда вычисления по формуле (15.118) на 6-разрядном десятичном компьютере дают значения $\psi(\alpha_0) = z(2, 0) \approx 2.69536 \cdot 10^7$, $\psi(\alpha_1) = z(2, -1) \approx -2.15629 \cdot 10^8$. В соответствии с формулой (15.118) получается следующее значение пристрелочного параметра:

$$\alpha = -1 - \frac{-2.15629 \cdot 10^8(-1 - 0)}{-2.15629 \cdot 10^8 - 2.69536 \cdot 10^7} \approx \alpha^* = -1 + 0.888888 = -0.111112.$$

Подстановка в формулы (15.116), (15.117) значения $\alpha = \alpha^*$ приводит к приближенному решению

$$y^*(x) = c_1^* e^{10x} + c_2^* e^{-10x} - \frac{1}{99} e^x,$$

$$z^*(x) = 10c_1^* e^{10x} - 10c_2^* e^{-10x} - \frac{1}{99} e^x$$

Здесь

$$c_1^* = \frac{1}{20} \left(\frac{1}{9} + \alpha^* \right) \approx \frac{1}{20} (0.111111 - 0.111112) = -5 \cdot 10^{-8},$$

$$c_2^* = \frac{1}{20} \left(\frac{1}{11} - \alpha^* \right) \approx \frac{1}{20} (0.0909091 + 0.111112) = 0.0101011$$

Тогда $c_1 - c_1^* \approx 5 \cdot 10^{-8}$, $c_2 - c_2^* \approx -10^{-7}$. Это означает, что

$$y(x) - y^*(x) \approx 5 \cdot 10^{-8} e^{10x} - 10^{-7} e^{-10x},$$

$$z(x) - z^*(x) \approx 5 \cdot 10^{-7} e^{10x} + 10^{-6} e^{-10x}$$

Наличие в погрешности компоненты, пропорциональной e^{10x} , приводит к тому, что при $x = 2$ погрешности решения достигают следующих значений $y(2) - y^*(2) \approx 243$, $z(2) - z^*(2) \approx 243$. ▲

§ 15.6. Дополнительные замечания

1. В этой главе метод конечных разностей и метод конечных элементов рассматривались только лишь применительно к решению двухточечных краевых задач для обыкновенных дифференциальных уравнений второго порядка. Значительно более широкую область применения этих методов представляют собой различные задачи для дифференциальных уравнений в частных производных. Ограниченный объем книги не позволяет отразить здесь богатство существующих подходов и разнообразие используемых приемов. Тем, кто интересуется решением уравнений в частных производных с помощью метода конечных разностей, рекомендуем первоначально обратиться к учебникам [9, 10*, 51, 73, 88], а затем — к книгам [67, 85, 86].

Как доступное введение в метод конечных элементов, можно рекомендовать книги [2, 34, 70, 93]. В дальнейшем следует обратиться к [43, 44, 68, 90, 105].

2. Так как мы рассматривали краевые задачи только для обыкновенных дифференциальных уравнений, то тем самым фактически лишили себя возможности обсуждать достоинства и недостатки метода конечных

разностей и метода конечных элементов в их сравнении между собой. Ограничимся лишь констатацией того, что для решения дифференциальных уравнений существуют два мощных метода, каждый из которых не обладает, вообще говоря, безусловным преимуществом над другим. Тем не менее отметим, что нередко наиболее эффективными оказываются именно те приближенные методы, которые сочетают в себе достоинства обоих методов.

3. При математическом моделировании различных физических явлений часто приходится решать краевые задачи, в которых дифференциальные уравнения или краевые условия являются нелинейными. Примером может служить дифференциальное уравнение.

$$-\frac{d}{dx} \left(k(x, u) \frac{du}{dx} \right) = f(x, y),$$

описывающее установившееся распределение тепла в стержне, теплофизические характеристики которого зависят от температуры u . Для решения таких задач широко используются метод конечных разностей и метод конечных элементов. Возникающие здесь дискретные краевые задачи нелинейны и требуют для вычисления решений использования специальных итерационных методов.

4. В последние десятилетия было осознано, что решение проблемы численного решения дифференциальных уравнений основано на использовании специальных методов теории приближения функций. Однако глубокая связь между проблемой аппроксимации функций и проблемой решения дифференциальных уравнений осталась за рамками данной книги. Отметим лишь, что приближенное решение $u^N(x)$, полученное с помощью проекционно-разностной схемы (15.81), (15.82), представляет собой линейный сплайн.

Глава 16

ЧИСЛЕННОЕ РЕШЕНИЕ ИНТЕГРАЛЬНЫХ УРАВНЕНИЙ

Интегральные уравнения часто возникают в различных приложениях. В терминах интегральных уравнений формулируются многие задачи астрофизики, геологии, теплотехники, теории упругости, теории пластичности, гидродинамики, ядерной физики, биомеханики, медицины и др. Нередко к решению интегральных уравнений сводят другие математические задачи, например некоторые задачи математической физики.

§ 16.1. Понятие об интегральных уравнениях

Уравнения, в которых искомая функция входит под знак интеграла, принято называть интегральными уравнениями. Наиболее часто встречаются интегральные уравнения Фредгольма¹ первого рода

$$\int_a^b K(x, s)u(s) \, ds = f(x), \quad x \in [a, b], \quad (16.1)$$

интегральные уравнения Фредгольма второго рода

$$u(x) = \int_a^b K(x, s)u(s) \, ds + f(x), \quad x \in [a, b], \quad (16.2)$$

а также интегральные уравнения Вольтерра^{2 3} первого рода

$$\int_a^x K(x, s)u(s) \, ds = f(x), \quad x \in [a, b] \quad (16.3)$$

¹ Фредгольм Эрик Ивар (1866—1927) — шведский математик. Является одним из создателей теории линейных интегральных уравнений.

² Вольтерра Вито (1860—1940) — итальянский математик один из основоположников теории интегральных уравнений. Создатель математической теории биологических сообществ, аппаратом которой служат дифференциальные и интегро-дифференциальные уравнения.

³ По правилам русского языка следовало бы писать «интегральное уравнение Вольтерры». Мы пишем «интегральное уравнение Вольтерра», следуя традиции, принятой в русскоязычной литературе по теории интегральных уравнений, в ущерб правописанию.

и интегральные уравнения Вольтерра второго рода

$$u(x) = \int_a^x K(x, s)u(s) \, ds + f(x), \quad x \in [a, b] \quad (16.4)$$

Функции $K(x, s)$ (*ядро интегрального уравнения*) и $f(x)$ (*правая часть интегрального уравнения*) считаются заданными. Под решением интегрального уравнения понимается функция $u(x)$, подстановка которой в уравнение обращает его в тождество.

При исследовании интегральных уравнений делаются различные предположения о свойствах ядра K и правой части f . Мы в основном ограничимся рассмотрением наиболее простого случая, когда обе эти функции непрерывны, хотя во многих прикладных задачах возникает необходимость решать уравнения с разрывными ядрами и правыми частями. Кроме того, будем считать эти функции вещественнонезначными, хотя нередко встречаются задачи, в которых ядра и правые части являются комплекснозначными. Решение u также будем искать в классе непрерывных вещественнонезначных функций.

Обратим внимание на то, что искомая функция u входит в уравнения (16.1)–(16.4) линейным образом. Такие интегральные уравнения принято называть *линейными*. При решении многих задач приходится иметь дело и с нелинейными интегральными уравнениями — аналогами уравнений (16.1)–(16.4). Это нелинейные интегральные уравнения Фредгольма первого рода

$$\int_a^b K(x, s, u(s)) \, ds = f(x), \quad x \in [a, b]$$

и второго рода

$$u(x) = \int_a^b K(x, s, u(s)) \, ds + f(x), \quad x \in [a, b],$$

а также нелинейные интегральные уравнения Вольтерра первого рода

$$\int_a^x K(x, s, u(s)) \, ds = f(x), \quad x \in [a, b]$$

и второго рода

$$u(x) = \int_a^x K(x, s, u(s)) \, ds + f(x), \quad x \in [a, b].$$

§ 16.2. Примеры задач, приводящих к интегральным уравнениям

Приведем несколько примеров задач, приводящих к интегральным уравнениям. Первые три из них — это примеры интегральных уравнений, возникающих при решении научных и технических задач. Последние три показывают, как к интегральным уравнениям могут быть сведены некоторые математические задачи.

1. Интегральное уравнение Абеля. В 1823 году Н. Абель, рассматривая задачу о движении материальной точки в вертикальной плоскости Oxz под действием силы тяжести, пришел к интегральному уравнению

$$\int_0^z \frac{u(s)}{\sqrt{z-s}} ds = f(z). \quad (16.5)$$

Здесь неизвестная функция u описывает кривую $x = u(z)$, по которой должна скатываться материальная точка так, чтобы начав свое движение без начальной скорости, она опускалась на расстояние z по вертикали за заданное время $t = f(z)$.

Эта задача была исторически первой задачей, оформленной как интегральное уравнение¹. Уравнение Абеля (16.5) является интегральным уравнением Вольтерра первого рода и возникает при решении ряда других задач, например при определении потенциальной энергии по периоду колебаний или при восстановлении рассеивающего поля по эффективному сечению.

2. Интегральное уравнение переноса излучения. При описании распространения светового излучения в атмосферах звезд и планет широко используется интегральное уравнение Фредгольма второго рода

$$S(\tau) = \frac{\omega}{2} \int_0^{\tau^*} E_1(|\tau' - \tau|) S(\tau') d\tau' + f(\tau), \quad (16.6)$$

которое в астрофизике называется *интегральным уравнением переноса излучения*.

Искомой здесь является функция источника $S(\tau)$, аргумент τ имеет смысл оптической глубины, τ^* — полная оптическая глубина атмосферы, функция $f(\tau)$ характеризует плотность источников излучения. Постоянная ω называется показателем альбедо (по физическому смыслу $0 < \omega < 1$).

¹ Сам термин «интегральное уравнение» был введен значительно позднее (в 1888 году) немецким математиком Диюба-Реймоном.

Функция $E_1(t) = \int_0^t \frac{e^{-sx}}{x} dx$ — это специальная функция, называемая интегральной показательной функцией первого порядка.

3. Задача вычисления входного сигнала. Действие многих приборов, регистрирующих физические поля, может быть описано с помощью интегрального уравнения Вольтерра первого рода

$$\int_0^t K(t, s) u(s) ds = f(t).$$

Здесь t — время; $u(t)$ — сигнал, поступающий на вход прибора; $f(t)$ — выходной сигнал; $K(t, s)$ — известная функция, определяемая устройством прибора.

Требуется восстановить входной сигнал $u(t)$ по регистрируемому на выходе из прибора выходному сигналу $f(t)$.

4. Связь задачи Коши с нелинейным интегральным уравнением Вольтерра. Обратим внимание на то, что функция $y(t)$ является решением задачи Коши

$$y'(t) = f(t, y(t)), \quad (16.7)$$

$$y(t_0) = y_0 \quad (16.8)$$

тогда и только тогда, когда она является решением нелинейного интегрального уравнения Вольтерра второго рода

$$y(t) = \int_{t_0}^t f(s, y(s)) ds + y_0. \quad (16.9)$$

Действительно, интегрируя уравнение (16.7) от t_0 до t и учитывая начальное условие (16.8), мы убеждаемся в том, что каждое решение задачи Коши одновременно является решением интегрального уравнения (16.9). Обратно, дифференцируя (16.9), мы приходим к (16.7), а подстановка $t = t_0$ в уравнение (16.9) дает (16.8). Таким образом, всякое решение интегрального уравнения (16.9) одновременно является решением задачи Коши (16.7), (16.8).

5. Связь краевой задачи с интегральным уравнением Фредгольма. Рассмотрим краевую задачу

$$-\frac{d}{dx} \left(k(x) \frac{du}{dx} \right) + q(x)u = f(x), \quad a \leq x \leq b, \quad (16.10)$$

$$u(a) = u_a, \quad u(b) = u_b. \quad (16.11)$$

Здесь $k(x)$, $q(x)$, $f(x)$ — непрерывные функции, заданные на отрезке $[a, b]$, причем $k(x) \geq k_0 > 0$

Положим $r(x) = \int_a^x \frac{dy}{k(y)}$, $R = \int_a^b \frac{dy}{k(y)}$. Дважды интегрируя уравнение

(16.10) по x и учитывая краевые условия (16.11), приходим к интегральному уравнению Фредгольма второго рода

$$u(x) = - \int_a^b G(x, s) q(s) u(s) ds + F(x), \quad a \leq x \leq b, \quad (16.12)$$

где

$$G(x, s) = \begin{cases} \frac{(R - r(s)) r(x)}{R} & \text{при } a \leq x \leq s \leq b, \\ \frac{r(x)(R - r(s))}{R} & \text{при } a \leq s \leq x \leq b, \end{cases}$$

$$F(x) = \int_a^b G(x, s) f(s) ds + \frac{R - r(x)}{R} u_a + \frac{r(x)}{R} u_b.$$

Так что всякое решение краевой задачи (16.10), (16.11) одновременно является и решением интегрального уравнения (16.12).

В свою очередь, дважды дифференцируя интегральное уравнение (16.12), мы приходим к уравнению (16.10). Кроме того, подстановка в уравнение (16.12) значений $x = a$ и $x = b$ дает (16.11). Следовательно, всякое решение интегрального уравнения (16.12) одновременно является решением задачи (16.10), (16.11).

Таким образом, краевая задача (16.10), (16.11) и интегральное уравнение (16.12) эквивалентны.

6. Связь с задачей дифференцирования. Пусть $f(x)$ — непрерывно дифференцируемая функция, заданная на отрезке $[a, b]$. Требуется найти ее производную

$$u(x) = f'(x), \quad a \leq x \leq b. \quad (16.13)$$

Интегрируя равенство (16.13) от a до x , мы приходим к равенству

$$\int_a^x u(s) ds = f(x) - f(a), \quad a \leq x \leq b. \quad (16.14)$$

Обратно, дифференцирование равенства (16.14) дает (16.13). Таким образом, задача вычисления производной функции f эквивалентна

задаче решения интегрального уравнения Вольтерра первого рода (16.14).

Аналогично, если функция f является n раз непрерывно дифференцируемой, то задача вычисления ее производной порядка n

$$u(x) = f^{(n)}(x), \quad a \leq x \leq b$$

эквивалентна задаче решения интегрального уравнения

$$\int_a^x \frac{(x-s)^{n-1}}{(n-1)!} u(s) ds = f(x) - \sum_{k=0}^{n-1} \frac{f^{(k)}(a)}{k!} (x-a)^k, \quad a \leq x \leq b. \quad (16.15)$$

§ 16.3. Интегральные уравнения Фредгольма второго рода. Введение в теорию

1. Обозначения и определения. Прежде чем перейти к изложению некоторых результатов из теории интегральных уравнений Фредгольма второго рода

$$u(x) = \int_a^b K(x, s) u(s) ds + f(x), \quad x \in [a, b] \quad (16.16)$$

и методов их численного решения, дадим некоторые определения и введем ряд обозначений.

Будем считать, что ядро K является непрерывной функцией, заданной на квадрате $\Pi = \{(x, s) | a \leq x \leq b, a \leq s \leq b\}$. Правую часть уравнения f будем также предполагать непрерывной функцией, заданной на отрезке $[a, b]$. Под решением уравнения (16.16) будем понимать непрерывную на отрезке $[a, b]$ функцию $u(x)$, подстановка которой в уравнение обращает его в тождество.

Определим скалярное произведение функций f и g формулой

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

Будем говорить, что функции f и g ортогональны, если $(f, g) = 0$.

Определим нормы

$$\|f\|_{C[a, b]} = \max_{[a, b]} |f(x)|, \quad \|f\| = \sqrt{\int_a^b |f(x)|^2 dx}$$

и заметим, что

$$\|f\| = \sqrt{(f, f)}.$$

Положим

$$\|K\|_C = \max_{(x, s) \in \Pi} |K(x, s)| \quad \text{и} \quad \|K\| = \sqrt{\int_a^b \int_a^b |K(x, s)|^2 dx ds}$$

Для частных производных функции K по переменным x и s будем использовать обозначения

$$K'_x = \frac{\partial}{\partial x} K, \quad K''_{xx} = \frac{\partial^2}{\partial x^2} K, \quad K'_s = \frac{\partial}{\partial s} K, \quad K''_{ss} = \frac{\partial^2}{\partial s^2} K.$$

Введем интегральный оператор \mathcal{K} и сопряженный к нему оператор \mathcal{K}^* :

$$\mathcal{K}u(x) = \int_a^b K(x, s)u(s) ds;$$

$$\mathcal{K}^*u(x) = \int_a^b K^*(x, s)u(s) ds,$$

где $K^*(x, s) = K(s, x)$.

В приложениях часто встречаются интегральные операторы с симметричными ядрами, то есть с ядрами, удовлетворяющими условию

$$K(x, s) \equiv K(s, x). \quad (16.17)$$

В этом случае $\mathcal{K}^* = \mathcal{K}$, то есть оператор \mathcal{K} является *самосопряженным*. Самосопряженные интегральные операторы обладают рядом замечательных свойств. Интегральные уравнения с такими операторами исследованы значительно более глубоко по сравнению с интегральными уравнениями общего вида, и они проще для численного решения.

Далее наряду с развернутой формой записи интегрального уравнения (16.16) будем использовать его краткую операторную форму записи

$$u = \mathcal{K}u + f.$$

2. Теория Фредгольма. Если правая часть $f(x) \equiv 0$, то уравнение (16.16) принимает вид

$$u(x) = \int_a^b K(x, s)u(s) ds, \quad x \in [a, b] \quad (16.18)$$

и называется *основным однородным уравнением*. Интегральное уравнение

$$u(x) = \int_a^b K^*(x, s)u(s) ds, \quad x \in [a, b] \quad (16.19)$$

называется *сопряженным однородным уравнением*.

Ясно, что решением каждого из уравнений (16.18) и (16.19) является функция $u(x) \equiv 0$, называемая *тривиальным решением*. Однако эти уравнения могут иметь и *нетривиальные решения* $u(x) \not\equiv 0$.

Основу теории интегрального уравнения (16.16) составляют следующие три теоремы.

Теорема 16.1 (первая теорема Фредгольма). Для того, чтобы при любой правой части f уравнение (16.16) имело единственное решение, необходимо и достаточно, чтобы основное однородное уравнение (16.18) имело лишь тривиальное решение. ■

Теорема 16.2 (вторая теорема Фредгольма). Основное однородное уравнение (16.18) и сопряженное однородное уравнение (16.19) имеют одинаковое и конечное число линейно независимых решений. ■

Теорема 16.3 (третья теорема Фредгольма). Для того, чтобы при заданной правой части f уравнение (16.16) было разрешимо, необходимо и достаточно, чтобы функция f была ортогональна всем решениям сопряженного однородного уравнения (16.19). ■

3. Корректность задачи. В случае, когда при любой правой части f уравнение

$$u = \mathcal{K}u + f \quad (16.20)$$

имеет единственное решение u , справедливы оценки

$$\|u\|_{C[a, b]} \leq C_1 \|f\|_{C[a, b]}, \quad (16.21)$$

$$\|u\| \leq C_2 \|f\| \quad (16.22)$$

с некоторыми постоянными C_1 и C_2 . В рассматриваемом случае задача решения интегрального уравнения Фредгольма второго рода (16.20) является корректной¹.

Действительно, пусть \tilde{u} — решение уравнения

$$\tilde{u} = \mathcal{K}\tilde{u} + \tilde{f} \quad (16.23)$$

с правой частью $\tilde{f} \approx f$. Вычитая друг из друга уравнения (16.20) и (16.23), убеждаемся, что погрешность $u - \tilde{u}$ удовлетворяет уравнению

$$u - \tilde{u} = \mathcal{K}(u - \tilde{u}) + f - \tilde{f}$$

¹ Напомним (см. § 3.1), что вычислительная задача называется корректной, если ее решение существует, единственно и устойчиво к малым погрешностям в данных.

Из неравенств (16.21), (16.22) немедленно следуют оценки

$$\|u - \tilde{u}\|_{C[a, b]} \leq C_1 \|f - \tilde{f}\|_{C[a, b]},$$

$$\|u - \tilde{u}\| \leq C_2 \|f - \tilde{f}\|,$$

которые говорят о непрерывной зависимости решения u от правой части f . Постоянные C_1 и C_2 играют здесь роль чисел обусловленности.

Приведем простые достаточные условия на ядро K , гарантирующие существование решения, его единственность и справедливость оценок (16.21), (16.22).

Теорема 16.4. Пусть ядро K удовлетворяет условию

$$q = \max_{x \in [a, b]} \int_a^b |K(x, s)| ds < 1. \quad (16.24)$$

Тогда при любой правой части f уравнение (16.20) имеет единственное решение u , причем справедлива оценка (16.21) с постоянной $C_1 = \frac{1}{1-q}$.

Доказательство. Пусть u — решение уравнения (16.20). Тогда

$$\begin{aligned} |u(x)| &= \left| \int_a^b K(x, s) u(s) ds + f(x) \right| \leq \int_a^b |K(x, s)| ds \|u\|_{C[a, b]} + \|f\|_{C[a, b]} \leq \\ &\leq q \|u\|_{C[a, b]} + \|f\|_{C[a, b]}, \quad x \in [a, b]. \end{aligned}$$

Отсюда

$$\|u\|_{C[a, b]} \leq q \|u\|_{C[a, b]} + \|f\|_{C[a, b]}.$$

Как следствие, верна оценка (16.21) с постоянной $C_1 = \frac{1}{1-q}$.

Пусть теперь u — решение уравнения (16.18), то есть решение уравнения (16.20) с правой частью $f(x) \equiv 0$. Тогда из оценки (16.21) следует, что $u(x) \equiv 0$. Поскольку основное однородное уравнение (16.18) имеет лишь тривиальное решение, то в силу теоремы 16.1 при любой правой части f уравнение (16.20) имеет единственное решение. ■

Аналогичным образом доказывается следующий результат.

Теорема 16.5. Пусть ядро K удовлетворяет условию

$$\|K\| = \sqrt{\int_a^b \int_a^b |K(x, s)|^2 dx ds} < 1. \quad (16.25)$$

Тогда при любой правой части f уравнение (16.20) имеет единственное решение, причем справедлива оценка (16.22) с постоянной $C_2 = \frac{1}{1 - \|K\|}$. ■

4. Связь невязки и погрешности. Пусть теперь u — точное решение интегрального уравнения (16.20), а \tilde{u} — некоторое приближенное решение того же уравнения. Определим *невязку*

$$r[\tilde{u}] = \tilde{u} - \mathcal{K}\tilde{u} - f, \quad (16.26)$$

отвечающую этому приближенному решению. В развернутой форме записи невязка задается формулой

$$r[\tilde{u}](x) = \tilde{u}(x) - \int_a^b K(x, s)\tilde{u}(s) ds - f(x), \quad x \in [a, b].$$

Заметим, что точному решению u отвечает нулевая невязка:

$$r[u] = u - \mathcal{K}u - f = 0. \quad (16.27)$$

Вычитая из равенства (16.27) равенство (16.26), убеждаемся в том, что погрешность $u - \tilde{u}$ является решением уравнения того же вида, что и (16.20), но с невязкой $r[\tilde{u}]$ в роли правой части f :

$$u - \tilde{u} = \mathcal{K}(u - \tilde{u}) + r[\tilde{u}].$$

Если справедливы неравенства (16.21) и (16.22), то погрешность $u - \tilde{u}$ можно оценить через невязку следующим образом:

$$\|u - \tilde{u}\|_{C[a, b]} \leq C_1 \|r[\tilde{u}]\|_{C[a, b]}, \quad (16.28)$$

$$\|u - \tilde{u}\| \leq C_2 \|r[\tilde{u}]\|. \quad (16.29)$$

Таким образом, если невязка, отвечающая приближенному решению \tilde{u} , достаточно мала, то малой будет и погрешность $u - \tilde{u}$.

5. Интегральное уравнение с приближенно заданным ядром. Предположим, что вместо ядра K интегрального уравнения

$$u(x) = \int_a^b K(x, s)u(s) ds + f(x), \quad x \in [a, b] \quad (16.30)$$

дано его приближение \tilde{K} , причем оба ядра заданы и непрерывны на квадрате Π . Пусть ядро K удовлетворяет условию (16.24), а погрешность $K - \tilde{K}$, с которой задано ядро, такова, что

$$\max_{(x \in [a, b])} \int_a^b |K(x, s) - \tilde{K}(x, s)| ds < \varepsilon,$$

где величина $\varepsilon > 0$ мала настолько, что $q + \varepsilon \leq \tilde{q} < 1$. В этом случае ядро \tilde{K} обладает свойством

$$\max_{[a, b]} \int_a^b |\tilde{K}(x, s)| ds \leq \tilde{q} < 1.$$

Поэтому в силу теоремы 16.4 решение \tilde{u} интегрального уравнения

$$\tilde{u}(x) = \int_a^b \tilde{K}(x, s) \tilde{u}(s) ds + f(x), \quad x \in [a, b] \quad (16.31)$$

с приближенно заданным ядром существует и единственno.

Вычитая из уравнения (16.30) уравнение (16.31), имеем

$$u(x) - \tilde{u}(x) = \int_a^b [K(x, s) - \tilde{K}(x, s)] \tilde{u}(s) ds + \int_a^b K(x, s) [u(s) - \tilde{u}(s)] ds.$$

Отсюда

$$\begin{aligned} |u(x) - \tilde{u}(x)| &\leq \int_a^b |K(x, s) - \tilde{K}(x, s)| ds \cdot \|\tilde{u}\|_{C[a, b]} + \\ &+ \int_a^b |K(x, s)| ds \cdot \|u - \tilde{u}\|_{C[a, b]}. \end{aligned}$$

Как следствие,

$$\|u - \tilde{u}\|_{C[a, b]} \leq \varepsilon \|\tilde{u}\|_{C[a, b]} + q \|u - \tilde{u}\|_{C[a, b]}.$$

В результате получаем неравенство

$$\frac{\|u - \tilde{u}\|_{C[a, b]}}{\|\tilde{u}\|_{C[a, b]}} \leq C_1 \varepsilon. \quad (16.32)$$

Оно говорит о том, что если ядро задано с достаточно малой погрешностью, то малой будет и относительная погрешность решения \tilde{u} .

Аналогичным образом обстоит дело тогда, когда выполнено условие (16.25), а погрешность $K - \tilde{K}$, с которой задано ядро, такова, что

$$\|K - \tilde{K}\| = \sqrt{\int_a^b \int_a^b |K(x, s) - \tilde{K}(x, s)|^2 dx ds} < \varepsilon,$$

где величина $\varepsilon > 0$ мала настолько, что $\|K\| + \varepsilon < 1$.

В этом случае решение уравнения (16.31) существует, единственно и справедлив аналог оценки (16.32) — оценка

$$\frac{\|u - \tilde{u}\|}{\|\tilde{u}\|} \leq C_2 \varepsilon. \quad (16.33)$$

6. Гладкость решений. Обратим сначала внимание на следующее стягивающее свойство оператора \mathcal{K} .

Лемма 16.1. Пусть ядро K непрерывно на квадрате Π и имеет непрерывные на Π частные производные $\frac{\partial K}{\partial x}, \frac{\partial^2 K}{\partial x^2}, \dots, \frac{\partial^n K}{\partial x^n}$. Если функция u интегрируема по Риману на отрезке $[a, b]$, то функция

$$z(x) = (\mathcal{K}u)(x) = \int_a^b K(x, s)u(s) ds$$

является n раз непрерывно дифференцируемой на $[a, b]$, причем

$$\|z^{(k)}\|_{C[a, b]} \leq M_k \sup_{a \leq x \leq b} |u(x)|, \quad 1 \leq k \leq n, \quad (16.34)$$

$$\text{где } M_k = \max_{x \in [a, b]} \left| \int_a^b \frac{\partial^k K}{\partial x^k}(x, s) ds \right|.$$

Доказательство. Как нетрудно видеть, справедливо равенство

$$z^{(k)}(x) = \int_a^b \frac{\partial^k K}{\partial x^k}(x, s)u(s) ds, \quad 1 \leq k \leq n,$$

из которого следует оценка (16.34). ■

Степень гладкости решения интегрального уравнения (16.30) определяется гладкостью ядра K и правой части f .

Теорема 16.6. Пусть ядро K удовлетворяет условиям леммы 16.1, а правая часть f является n раз непрерывно дифференцируемой на отрезке $[a, b]$ функцией. Тогда решение интегрального уравнения

(16.30) также является n раз непрерывно дифференцируемой на отрезке $[a, b]$ функцией, причем

$$\|u^{(k)}\|_{C[a, b]} \leq M_k \|u\|_{C[a, b]} + \|f^{(k)}\|_{C[a, b]}, \quad 1 \leq k \leq n, \quad (16.35)$$

где постоянные M_k — те же, что и в лемме 16.1

Доказательство. В силу структуры уравнения (16.30) его решение представляется в виде суммы

$$u = z + f, \quad (16.36)$$

где $z = \mathcal{K}u$. Поэтому

$$u^{(k)} = z^{(k)} + f^{(k)}, \quad 1 \leq k \leq n.$$

Как следствие,

$$\|u^{(k)}\|_{C[a, b]} \leq \|z^{(k)}\|_{C[a, b]} + \|f^{(k)}\|_{C[a, b]}, \quad 1 \leq k \leq n.$$

Для доказательства оценки (16.35) осталось воспользоваться неравенством (16.34) и тем, что для непрерывной на $[a, b]$ функции u верно равенство $\sup_{a \leq x \leq b} |u(x)| = \|u\|_{C[a, b]}$. ■

Как мы увидим ниже, точность численных методов решения интегрального уравнения существенно зависит от степени гладкости решения u . Чем более гладким является решение, тем с более высокой точностью его можно вычислить.

Если ядро K и правая часть f достаточно гладкие, то в силу теоремы 16.6 достаточно гладким будет и решение. Однако на практике часто встречаются задачи, в которых ядро является гладким, а правая часть f не является гладкой и может даже быть разрывной функцией.

Рассмотрим ситуацию, когда ядро K удовлетворяет условиям леммы 16.1, а про правую часть f известно, что она лишь интегрируема по Риману на $[a, b]$. Тогда решение уравнения в силу формулы (16.36) будет лишь интегрируемым по Риману и не будет гладким. Более того, каждая точка разрыва функции f будет и точкой разрыва функции u .

В этой ситуации возможен следующий выход из положения. Подставив формулу $u = z + f$ в уравнение $u = \mathcal{K}u + f$, убеждаемся, что слагаемое $z = \mathcal{K}u$ является решением интегрального уравнения того же вида, что и исходное:

$$z = \mathcal{K}z + F, \quad (16.37)$$

но с правой частью $F = \mathcal{K}f$ в роли f . Решение этого уравнения уже является n раз непрерывно дифференцируемой функцией. Вычислив приближение \tilde{z} к z с погрешностью $\varepsilon = \tilde{z} - z$, мы можем затем найти приближение $\tilde{u} = \tilde{z} + f$ к u . Оно будет иметь ту же погрешность $u - \tilde{u} = \varepsilon$, поскольку $u - \tilde{u} = z - \tilde{z}$.

Конечно, следует иметь в виду, что платой за такой выход из положения является необходимость в соответствующих методах вычислять значения функции $F(x) = \int_a^b K(x, s) f(s) ds$ вместо значений функции f .

§ 16.4. Интегральные уравнения Фредгольма второго рода. Метод квадратур

1. Метод квадратур. Одним из наиболее простых и распространенных методов решения интегрального уравнения Фредгольма второго рода

$$u(x) = \int_a^b K(x, s) u(s) ds + f(x), \quad x \in [a, b] \quad (16.38)$$

с гладкими ядрами является *метод квадратур*, основанный на аппроксимации значений \mathcal{K} интегрального оператора (т.е. на приближенной замене интеграла, входящего в уравнение (16.38)) с использованием одной из квадратурных формул (см. гл. 13).

Пусть за основу взята квадратурная формула

$$\int_a^b g(s) ds \approx \sum_{j=0}^N c_j g(x_j) \quad (16.39)$$

с узлами x_j (где $a \leq x_0 < x_1 < \dots < x_N \leq b$) и весами $c_j > 0$.

Используя формулу (16.39) с $g(s) = K(x, s) u(s)$ для приближенного вычисления значения интегрального оператора, имеем

$$\int_a^b K(x, s) u(s) ds \approx \sum_{j=0}^N c_j K(x, x_j) u(x_j),$$

и от уравнения (16.38) мы приходим к приближенному равенству

$$u(x) = \sum_{j=0}^N c_j K(x, x_j) u(x_j) + f(x), \quad x \in [a, b].$$

Функцию u^N , удовлетворяющую уравнению

$$u^N(x) = \sum_{j=0}^N c_j K(x, x_j) u^N(x_j) + f(x), \quad x \in [a, b] \quad (16.40)$$

примем за приближенное решение уравнения (16.38).

Если уравнение (16.40) имеет решение, то его можно найти следующим образом. Подставляя в (16.40) значения $x = x_i$, для $i = 0, 1, \dots, N$, приходим к системе линейных алгебраических уравнений

$$u_i = \sum_{j=0}^N c_j K_{ij} u_j + f_i, \quad i = 0, 1, \dots, N \quad (16.41)$$

относительно вектора неизвестных $\mathbf{u}_N = (u_0, u_1, \dots, u_N)^T$, где $u_j = u^N(x_j)$. Здесь и ниже используются обозначения $K_{ij} = K(x_i, x_j)$ и $f_i = f(x_i)$.

Замечание. Обратим внимание на то, что решение уравнения (16.40) однозначно определяется своими значениями $u_j = u^N(x_j)$ в узлах x_j , $0 \leq j \leq N$. Поэтому уравнение (16.40) однозначно разрешимо тогда и только тогда, когда однозначно разрешима система (16.41).

Если ввести матрицу \mathbf{B}_N с элементами $b_{ij} = c_j K_{ij}$, $0 \leq i, j \leq N$ и вектор $\mathbf{f}_N = (f_0, f_1, \dots, f_N)^T$, то систему (16.41) можно записать в матричном виде:

$$\mathbf{u}_N = \mathbf{B}_N \mathbf{u}_N + \mathbf{f}_N. \quad (16.42)$$

Эту же систему можно записать в виде

$$\mathbf{A}_N \mathbf{u}_N = \mathbf{f}_N, \quad (16.43)$$

где $\mathbf{A}_N = \mathbf{E}_N - \mathbf{B}_N$, а \mathbf{E}_N — единичная матрица.

При умеренном числе N неизвестных решение системы уравнений (16.43) можно найти с помощью одного из прямых методов (см. гл. 5). При очень большом числе неизвестных следует использовать один из итерационных методов (см. гл. 6), например, — метод простой итерации

$$\mathbf{u}_N^{(k+1)} = \mathbf{B}_N \mathbf{u}_N^{(k)} + \mathbf{f}_N, \quad k \geq 0. \quad (16.44)$$

Заметим, что

$$\sum_{j=0}^N c_j |K_{ij}| \approx \int_a^b |K(x_i, s)| \, ds.$$

Поэтому, если выполнено предположение (16.24), то можно рассчитывать на выполнение условия

$$\|\mathbf{B}_N\|_\infty = \max_{0 \leq i \leq N} \sum_{j=0}^N c_j |K_{ij}| < 1,$$

гарантирующего не только существование и единственность решения системы (16.42), но и сходимость метода простой итерации (16.44) (см. § 6.1).

Как мы видели раньше, системы уравнений с симметричными матрицами предпочтительнее для численного решения, чем системы

с несимметричными матрицами. Обратим внимание на то, что даже если решается интегральное уравнение с симметричным ядром K , матрица системы (16.43) не обязана быть симметричной. Она будет симметричной, если $c_1 = c_2 = \dots = c_N$.

Если ядро симметрично, то систему (16.41) можно симметризовать, умножив i -е уравнение системы на c_i . В результате мы получим систему

$$c_i u_i = \sum_{j=0}^N c_i K_{ij} c_j u_j + c_i f_i, \quad i = 0, 1, \dots, N$$

с симметричной матрицей.

Другой способ симметризации состоит в умножении i -го уравнения системы (16.41) на $\sqrt{c_i}$. Сделав замену $v_i = \sqrt{c_i} u_i$, мы приходим к системе уравнений

$$v_i = \sum_{j=0}^N \sqrt{c_i} K_{ij} \sqrt{c_j} v_j + \sqrt{c_i} f_i, \quad i = 0, 1, \dots, N$$

с симметричной матрицей.

Тем или иным образом решив систему (16.43), мы затем можем вычислить приближенное решение u^N во всех точках отрезка $[a, b]$ по формуле (16.40).

Замечание. Приближенные значения решения в точках отрезка $[a, b]$, отличных от узловых, можно восстановить и с помощью одного из методов интерполяции (см. гл. 11). Тем не менее формула (16.40) выглядит предпочтительнее.

Покажем, как для аппроксимации интегрального уравнения (16.38) могут быть использованы простейшие квадратурные формулы.

2. Применение формулы прямоугольников. Пусть $a = x_0 < x_1 < \dots < x_N = b$, $x_{j-1/2} = (x_{j-1} + x_j)/2$, $h_j = x_j - x_{j-1}$, $1 \leq j \leq N$. Напомним, что составная квадратурная формула прямоугольников имеет вид

$$\int_a^b g(x) dx \approx \sum_{j=1}^N g(x_{j-1/2}) h_j. \quad (16.45)$$

Применяя эту формулу для аппроксимации значений интегрального оператора, приходим к уравнению

$$u^N(x) = \sum_{j=1}^N h_j K(x, x_{j-1/2}) u^N(x_{j-1/2}) + f(x), \quad x \in [a, b] \quad (16.46)$$

для приближенного решения u^N . Подставляя в это уравнение значения $x = x_{i-1/2}$ для $i = 1, 2, \dots, N$, выводим систему линейных алгебраических уравнений относительно значений $u_{i-1/2} = u^N(x_{i-1/2})$, $1 \leq i \leq N$:

$$u_{i-1/2} = \sum_{j=1}^N h_j K_{i-1/2, j-1/2} u_{j-1/2} + f_{i-1/2}, \quad i = 1, 2, \dots, N. \quad (16.47)$$

Здесь $K_{i-1/2, j-1/2} = K(x_{i-1/2}, x_{j-1/2})$; $f_{i-1/2} = f(x_{i-1/2})$.

Систему (16.47) можно записать в матричном виде

$$\mathbf{u}_N = \mathbf{B}_N \mathbf{u}_N + \mathbf{f}_N.$$

Здесь $\mathbf{u}_N = (u_{1/2}, u_{3/2}, \dots, u_{N-1/2})^T$ — вектор неизвестных; \mathbf{B}_N — квадратная матрица с элементами $b_{ij} = h_j K_{i-1/2, j-1/2}$, $1 \leq i, j \leq N$, а $\mathbf{f}_N = (f_{1/2}, f_{3/2}, \dots, f_{N-1/2})^T$.

После того, как эта система решена, приближенное решение может быть найдено во всех точках отрезка $[a, b]$ подстановкой найденных значений $u_{i-1/2} = u^N(x_{i-1/2})$ в правую часть формулы (16.46).

Приведем две теоремы, дающие теоретическое обоснование метода квадратур с использованием формулы прямоугольников (16.45).

Теорема 16.7. Пусть ядро K непрерывно на квадрате Π и имеет непрерывную на Π частную производную K'_s . Пусть также выполнены условие (16.24) и условие

$$h_{\max} \leq \frac{1-q}{2M}, \quad (16.49)$$

$$\text{где } M = \max_{a \leq x \leq b} \int_a^b |K'_s(x, s)| \, ds.$$

Тогда решение уравнения (16.46) существует, единственно и для него справедлива оценка

$$\|u^N\|_{C[a, b]} \leq \frac{2}{1-q} \|f\|_{C[a, b]}. \quad (16.50)$$

Доказательство. Покажем сначала, что справедливо неравенство

$$\max_{a \leq x \leq b} \sum_{j=1}^N h_j |K(x, x_{j-1/2})| \leq \frac{1+q}{2} < 1. \quad (16.51)$$

Положим

$$\varepsilon(x) = \int_a^b |K(x, s)| \, ds - \sum_{j=1}^N h_j |K(x, x_{j-1/2})|$$

и заметим, что

$$|\varepsilon(x)| = \left| \sum_{j=1}^N \int_{x_{j-1}}^{x_j} [|K(x, s) - K(x, x_{j-1/2})|] ds \right| \leq \sum_{j=1}^N \int_{x_{j-1}}^{x_j} |K(x, s) - K(x, x_{j-1/2})| ds .$$

Пользуясь для $s \in [x_{j-1}, x_j]$ неравенством

$$|K(x, s) - K(x, x_{j-1/2})| = \left| \int_{x_{j-1/2}}^s K'_s(x, t) dt \right| \leq \int_{x_{j-1}}^{x_j} |K'_s(x, s)| ds$$

и предположением (16.49), приходим к неравенству

$$|\varepsilon(x)| \leq \sum_{j=1}^N \int_{x_{j-1}}^{x_j} |K'_s(x, s)| ds h_j \leq \int_a^b |K'_s(x, s)| ds h_{\max} \leq M h_{\max} \leq \frac{1-q}{2} .$$

Таким образом,

$$\sum_{j=1}^N h_j |K(x, x_{j-1/2})| \leq \int_a^b |K(x, s)| ds + |\varepsilon(x)| \leq q + \frac{1-q}{2} = \frac{1+q}{2} < 1 .$$

Неравенство (16.51) доказано.

Из него следует, что

$$\|\mathbf{B}_N\|_\infty = \max_{1 \leq i \leq N} \sum_{j=1}^N |b_{ij}| = \max_{1 \leq i \leq N} \sum_{j=1}^N h_j |K(x_{i-1/2}, x_{j-1/2})| \leq \frac{1+q}{2} < 1 .$$

Поскольку $\|\mathbf{B}_N\|_\infty < 1$, то решение системы (16.48) существует и единственno, а следовательно существует и единственno решение u^N уравнения (16.46).

Заметим, что из (16.46) следует неравенство

$$|u^N(x)| \leq \sum_{j=1}^N h_j |K(x, x_{j-1/2})| \|u^N\|_{C[a, b]} + \|f\|_{C[a, b]} .$$

Откуда, используя неравенство (16.51), выводим

$$\|u^N\|_{C[a, b]} \leq \frac{1+q}{2} \|u^N\|_{C[a, b]} + \|f\|_{C[a, b]} .$$

Таким образом,

$$\|u^N\|_{C[a,b]} \leq \frac{1}{1-(1+q)/2} \|f\|_{C[a,b]} = \frac{2}{1-q} \|f\|_{C[a,b]}. \blacksquare$$

Замечание. Поскольку при выполнении условия (16.49) справедливо неравенство $\|B_N\|_\infty < 1$, то для решения системы (16.48) можно использовать метод простой итерации

$$u_N^{(k+1)} = B_N u_N^{(k)} + f_N, \quad k \geq 0$$

с произвольным начальным приближением $u_N^{(0)}$.

Теорема 16.8. Пусть ядро K дважды непрерывно дифференцируемо на квадрате Π , а функция f дважды непрерывно дифференцируема на отрезке $[a, b]$. Пусть также выполнены условия (16.24) и (16.49).

Тогда справедлива оценка погрешности

$$\|u - u^N\|_{C[a,b]} \leq Ch_{\max}^2, \quad (16.52)$$

$$\text{где } C = \frac{b-a}{12(1-q)} \max_{(x,s) \in \Pi} \left| \frac{\partial^2}{\partial s^2} (K(x,s)u(s)) \right|.$$

Доказательство. В силу теоремы 16.6 функция u является дважды непрерывно дифференцируемой на отрезке $[a, b]$.

Пусть

$$R^N(x) = \int_a^b K(x,s)u(s) ds - \sum_{j=1}^N h_j K(x, x_{j-1/2})u(x_{j-1/2}).$$

В силу теоремы об оценке погрешности составной формулы прямоугольников справедливо неравенство

$$|R^N(x)| \leq \frac{b-a}{24} \max_{a \leq s \leq b} \left| \frac{\partial^2}{\partial s^2} (K(x,s)u(s)) \right| h_{\max}^2. \quad (16.53)$$

Вычитая из уравнения (16.38) уравнение (16.46), имеем

$$u(x) - u^N(x) = \sum_{j=1}^N h_j K(x, x_{j-1/2})[u(x_{j-1/2}) - u^N(x_{j-1/2})] + R^N(x).$$

Таким образом, погрешность $u - u^N$ является решением уравнения вида (16.46) с R^N в роли f . Используя неравенства (16.50) и (16.53), приходим к оценке

$$\begin{aligned} \|u - u^N\|_{C[a,b]} &\leq \frac{2}{1-q} \|R^N\|_{C[a,b]} \leq \\ &\leq \frac{2}{1-q} \frac{b-a}{24} \max_{(x,s) \in \Pi} \left| \frac{\partial^2}{\partial s^2} (K(x,s)u(s)) \right| h_{\max}^2 = Ch_{\max}^2. \blacksquare \end{aligned}$$

3. Применение формулы трапеций. Возьмем за основу составную квадратурную формулу трапеций

$$\int_a^b g(x) dx \approx g(x_0) \frac{h_1}{2} + \sum_{j=1}^{N-1} g(x_j) h_{j+1/2} + g(x_N) \frac{h_N}{2},$$

где $h_{j+1/2} = (h_j + h_{j+1})/2$. Применяя эту формулу для аппроксимации значений интегрального оператора, приходим к уравнению

$$\begin{aligned} u^N(x) &= K(x, x_0)u^N(x_0) \frac{h_1}{2} + \sum_{j=1}^{N-1} K(x, x_j)u^N(x_j)h_{j+1/2} + \\ &+ K(x, x_N)u^N(x_N) \frac{h_N}{2} + f(x), \quad x \in [a, b] \end{aligned} \quad (16.54)$$

для приближенного решения u^N . Подставляя в это уравнения значения $x = x_i$ для $i = 0, 1, \dots, N$, выводим систему линейных алгебраических уравнений относительно значений $u_i = u^N(x_i)$:

$$u_i = K_{i0} \frac{h_1}{2} u_0 + \sum_{j=1}^{N-1} K_{ij} h_{j+1/2} u_j + K_{iN} \frac{h_N}{2} u_N + f_i, \quad 0 \leq i \leq N. \quad (16.55)$$

Решив эту систему, мы можем затем найти приближенное решение во всех точках отрезка $[a, b]$, используя формулу (16.54).

Следующие две теоремы доказываются аналогично теоремам 16.7 и 16.8.

Теорема 16.9. Пусть ядро K непрерывно на квадрате Π и имеет непрерывную на Π частную производную K'_s . Пусть также выполнены условия (16.24) и (16.49). Тогда решение уравнения (16.54) существует, единственно и для него справедлива оценка (16.50). ■

Теорема 16.10. Пусть ядро K дважды непрерывно дифференцируемо на квадрате Π , а функция f дважды непрерывно дифференцируема на отрезке $[a, b]$. Пусть также выполнены условия (16.24) и (16.49).

Тогда справедлива оценка погрешности

$$\|u - u^N\|_{C[a, b]} \leq 2Ch_{\max}^2, \quad (16.56)$$

где постоянная C та же, что и в (16.52). ■

Замечание. При выполнении условия (16.49) решение системы (16.55) может быть найдено с помощью метода простой итерации

$$u_i^{(k+1)} = K_{10} \frac{h_1}{2} u_0^{(k)} + \sum_{j=1}^{N-1} K_{1j} h_{j+1/2} u_j^{(k)} + K_{IN} \frac{h_N}{2} u_N^{(k)} + f_i, \quad 0 \leq i \leq N$$

при любом начальном приближении $u_N^{(0)} = (u_0^{(0)}, u_1^{(0)}, \dots, u_N^{(0)})^T$.

4. Погрешность метода квадратур. Обратим внимание на то, что сделанные в теоремах 16.8 и 16.10 предположения о гладкости ядра K и правой части f были нужны лишь для доказательства оценок (16.52) и (16.56), говорящих о том, что метод квадратур с использованием формул прямоугольников и трапеций имеет второй порядок точности. Можно показать, что и в случае, когда ядро K и правая часть f являются только непрерывными, при выполнении условия (16.24) метод квадратур является сходящимся:

$$\|u - u^N\|_{C[a, b]} \rightarrow 0 \quad \text{при } h_{\max} \rightarrow 0.$$

Приведенные в теоремах 16.8 и 16.10 оценки погрешностей являются априорными. Они указывают на то, что рассматриваемые методы обладают вторым порядком точности. Однако они непригодны для практической оценки погрешностей полученных приближений. Для того чтобы оценить погрешность, нужно использовать какие-либо апостериорные оценки.

Аналогично тому, как это делалось для оценки погрешности численного интегрирования, можно применять правило Рунге. Пусть используются методы на основе формул прямоугольников или трапеций с постоянным шагом $h = (b - a)/N$, где N — четное число. Тогда, вычислив приближенные решения u^N с шагом h и $u^{N/2}$ с шагом $2h$, можно оценить погрешность приближения u^N по формуле

$$u(x) - u^N(x) \approx \frac{u^N(x) - u^{N/2}(x)}{3}.$$

Распространенным приемом, применяемым для оценки качества приближенного решения является вычисление невязки

$$r[u^N] = u^N - \mathcal{K}u^N - f.$$

Оценки (16.28) и (16.29), примененные к $\tilde{u} = u^N$, дают неравенства

$$\|u - u^N\|_{C[a, b]} \leq C_1 \|r[u^N]\|_{C[a, b]},$$

$$\|u - u^N\| \leq C_2 \|r[u^N]\|,$$

позволяющие судить о величине погрешности $u - u^N$ по величине невязки $r[u^N]$.

получать методы, формально обладающие более высоким порядком точности. Например, формула Симпсона и формула Милна (Боде) (см. § 13.2) приведут к методам, обладающим четвертым и шестым порядком точности соответственно. Популярным подходом к построению методов решения интегральных уравнений является использование квадратурных формул Гаусса (см. § 13.3). Однако следует иметь в виду, что высокая точность этих формул предполагает, что решение u и ядро K являются достаточно гладкими функциями. Если это не так, то переход к использованию формул высокого порядка точности, скорее всего, не даст желаемого результата.

В случае, когда ядро K интегрального уравнения является достаточно гладким, а правая часть f не является гладкой или даже разрывна, можно использовать описанный в п. 6 § 16.2 прием перехода к представлению решения u в виде $u = z + f$, где z является решением интегрального уравнения $z = \mathcal{K}z + F$ с гладкой правой частью $F = \mathcal{K}u$. Вычислив методом квадратур приближение z^N к решению этого уравнения, мы можем затем вычислить приближение $u^N = z^N + f$ к решению u .

§ 16.5. Интегральные уравнения Фредгольма второго рода. Проекционные методы

Современной альтернативой методу квадратур численного решения интегрального уравнения Фредгольма второго рода

$$u = \mathcal{K}u + f \quad (16.57)$$

является класс методов, которые принято называть *проекционными*. Опишем некоторые из этих методов и начнем с классического варианта метода Галеркина (в теории интегральных уравнений этот метод часто называют *методом моментов*).

1. Метод Галеркина. Будем искать приближенное решение уравнения (16.57) в виде линейной комбинации заданных линейно независимых базисных функций $\varphi_1(x), \varphi_2(x), \dots, \varphi_N(x)$:

$$u^N(x) = \sum_{j=1}^N \alpha_j \varphi_j(x), \quad (16.58)$$

где коэффициенты $\alpha_1, \alpha_2, \dots, \alpha_N$ неизвестны и подлежат определению.

Заметим, что невязка, соответствующая приближенному решению u^N , имеет вид

$$r[u^N] = u^N - \mathcal{K}u^N - f = \sum_{j=1}^N \alpha_j \varphi_j - \sum_{j=1}^N \alpha_j \mathcal{K}\varphi_j - f. \quad (16.59)$$

Потребуем, чтобы невязка была ортогональна всем базисным функциям:

$$(r[u^N], \varphi_i) = 0, \quad 1 \leq i \leq N. \quad (16.60)$$

Подставляя в эти соотношения выражение (16.59), приходим к системе линейных алгебраических уравнений относительно вектора неизвестных $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_N)^T$:

$$\sum_{j=1}^N (\varphi_j, \varphi_i) \alpha_j = \sum_{j=1}^N (\mathcal{K}\varphi_j, \varphi_i) \alpha_j + (f, \varphi_i), \quad 1 \leq i \leq N$$

или, с использованием матричной формы записи, к системе вида

$$D_N \alpha = K_N \alpha + f_N. \quad (16.61)$$

Здесь D_N, K_N — квадратные матрицы с элементами

$$d_{ij} = (\varphi_j, \varphi_i) = \int_a^b \varphi_j(x) \varphi_i(x) dx, \quad 1 \leq i, j \leq N,$$

$$K_{ij} = (\mathcal{K}\varphi_j, \varphi_i) = \int_a^b \int_a^b K(x, s) \varphi_j(s) \varphi_i(x) dx ds, \quad 1 \leq i, j \leq N,$$

$f_N = (f_1, f_2, \dots, f_N)^T$ — вектор-столбец с элементами

$$f_i = (f, \varphi_i) = \int_a^b f(x) \varphi_i(x) dx, \quad 1 \leq i \leq N. \quad (16.62)$$

Замечание. Обратим внимание на то, что матрица D_N является симметричной. Если ядро интегрального оператора симметрично (т.е. если $K(x, s) \equiv K(s, x)$), то матрица K_N также является симметричной.

2. Метод Канторовича. Часто в методе Галеркина приближенное решение ищется не в виде (16.58), а в виде

$$u^N(x) = \sum_{j=1}^N \alpha_j \varphi_j(x) + f(x). \quad (16.63)$$

В этом случае невязка имеет вид

$$r[u^N] = u^N - \mathcal{K}u^N - f = \sum_{j=1}^N \alpha_j \varphi_j - \sum_{j=1}^N \alpha_j \mathcal{K}\varphi_j - \mathcal{K}f. \quad (16.64)$$

Она отличается от невязки (16.59) тем, что в ней правая часть f заменена на $F = \mathcal{K}f$. Поэтому система уравнений (16.60) принимает вид

$$\mathbf{D}_N \boldsymbol{\alpha} = \mathbf{K}_N \boldsymbol{\alpha} + \mathbf{F}_N$$

с теми же матрицами \mathbf{D}_N и \mathbf{K}_N , что и в (16.61), но с $\mathbf{F}_N = (F_1, F_2, \dots, F_N)^T$ вместо \mathbf{f}_N , где

$$F_i = (F, \varphi_i) = \int_a^b F(x) \varphi_i(x) dx = \int_a^b \int K(x, s) f(s) \varphi_i(x) dx ds, \quad 1 \leq i \leq N. \quad (16.65)$$

Описанную модификацию метода Галеркина иногда называют *методом Канторовича*¹.

Соображения, приводящие к тому, что приближенное решение имеет смысл искать в виде (16.63), были приведены в п. 6 § 16.2. Напомним, что этот подход имеет смысл, если ядро интегрального уравнения является гладким, а правая часть не является достаточно гладкой функцией. Тогда, представляя решение u в виде $u = z + f$, мы можем найти слагаемое z , решая интегральное уравнение

$$z = \mathcal{K}z + F \quad (16.66)$$

с гладкой правой частью $F = \mathcal{K}f$. Метод Канторовича можно рассматривать как классический метод Галеркина, примененный к решению уравнения (16.66).

Следует подчеркнуть, что цена, которую мы платим за возможность повышения точности приближения состоит в усложнении формулы (16.65) вычисления правой части системы по сравнению с формулой (16.62).

3. Итерированные методы Галеркина и Канторовича. После того, как методом Галеркина найдено приближенное решение u^N , можно найти уточненное приближение \tilde{u}^N , подставив u^N в правую часть интегрального уравнения:

$$\tilde{u}^N = \mathcal{K}u^N + f. \quad (16.67)$$

Этот метод принято называть *итерированным методом Галеркина* или *методом Слоана*.

Аналогично, после того как методом Канторовича найдено приближенное решение u^N , можно найти уточненное приближение по формуле (16.67). В этом состоит *итерированный метод Канторовича*.

¹ Леонид Витальевич Канторович (1912—1986) — советский математик и экономист, лауреат Нобелевской премии по экономике, один из создателей линейного программирования.

4. Проекционные методы с использованием кусочно-постоянных базисных функций. Наиболее просто описанные выше проекционные методы реализуются при использовании в качестве базисных функций кусочно-постоянных функций.

Введем на отрезке $[a, b]$ сетку с узлами $a = x_0 < x_1 < \dots < x_N = b$ и шагами $h_j = x_j - x_{j-1}$. В качестве базисных функций используем функции $\varphi_j = \chi_j$, где

$$\chi_j(x) = \begin{cases} 1 & \text{для } x \in [x_{j-1}, x_j), \\ 0 & \text{для } x \notin [x_{j-1}, x_j). \end{cases} \quad 1 \leq j \leq N,$$

$$\chi_N(x) = \begin{cases} 1 & \text{для } x \in [x_{N-1}, x_N], \\ 0 & \text{для } x \notin [x_{N-1}, x_N]. \end{cases}$$

В этом случае приближение (16.58) будет кусочно-постоянной функцией, принимающей значения α_j на промежутках $[x_{j-1}, x_j)$, $1 \leq j \leq N$. Имеет смысл обозначить эти значения через $u_{j-1/2}$, приближенное решение обозначить через u^h и записать его в виде

$$u^h(x) = \sum_{j=1}^N u_{j-1/2} \chi_j(x). \quad (16.68)$$

При таком выборе базисных функций система уравнений метода Галеркина (16.59) принимает вид

$$h_i u_{i-1/2} = \sum_{j=1}^N \int_{x_{j-1}}^{x_i} \left[\int_{x_{j-1}}^{x_j} K(x, s) ds \right] dx \cdot u_{j-1/2} + \int_{x_{i-1}}^{x_i} f(x) dx, \quad 1 \leq i \leq N.$$

Делением i -го уравнения системы на h_i ее можно преобразовать к виду

$$u_N = B_N u_N + f_N, \quad (16.69)$$

где $u_N = (u_{1/2}, \dots, u_{N-1/2})^T$ — вектор неизвестных; B_N — матрица с элементами

$$b_{ij} = \frac{1}{h_i} \int_{x_{i-1}}^{x_i} \left[\int_{x_{j-1}}^{x_j} K(x, s) ds \right] dx, \quad 1 \leq i, j \leq N, \text{ а } f_N = (f_{1/2}, \dots, f_{N-1/2})^T —$$

$$\text{вектор с элементами } f_{i-1/2} = \frac{1}{h_i} \int_{x_{i-1}}^{x_i} f(x) dx, \quad 1 \leq i \leq N.$$

Заметим, что

$$\begin{aligned}\|\mathbf{B}_N\|_{\infty} &= \max_{1 \leq i \leq N} \sum_{j=1}^N |b_{ij}| \leq \max_{1 \leq i \leq N} \frac{1}{h_i} \int_{x_{i-1}}^{x_i} \left[\sum_{j=1}^N \int_{x_{j-1}}^{x_j} |K(x, s)| ds \right] dx = \\ &= \max_{1 \leq i \leq N} \frac{1}{h_i} \int_{x_{i-1}}^{x_i} \left[\int_a^b |K(x, s)| ds \right] dx \leq q = \max_{a \leq x \leq b} \int_a^b |K(x, s)| ds.\end{aligned}$$

Таким образом, если выполнено условие (16.24), то $\|\mathbf{B}_N\|_{\infty} \leq q < 1$ и поэтому система (16.69) однозначно разрешима. Кроме того, ее решение может быть найдено методом простой итерации

$$\mathbf{u}_N^{(k+1)} = \mathbf{B}_N \mathbf{u}_N^{(k)} + \mathbf{f}_N, \quad k \geq 0$$

при произвольном начальном приближении $\mathbf{u}^{(0)}$.

После того как галерkinское приближение (16.68) найдено, итерированным методом Галеркина (16.67) можно вычислить приближение

$$\tilde{\mathbf{u}}^h(x) = \sum_{j=1}^N \int_{x_{j-1}}^{x_j} K(x, s) ds \cdot u_{j-1/2} + f(x), \quad x \in [a, b].$$

В методе Канторовича с использованием кусочно-постоянных базисных функций приближенное решение ищется в виде

$$\mathbf{u}^h(x) = \sum_{j=1}^N z_{j-1/2} \chi_j(x) + f(x), \quad (16.70)$$

Соответствующая система уравнений для определения вектора неизвестных $\mathbf{z}_N = (z_{1/2}, \dots, z_{N-1/2})^T$ имеет вид

$$\mathbf{z}_N = \mathbf{B}_N \mathbf{z}_N + \mathbf{F}_N, \quad (16.71)$$

где \mathbf{B}_N — та же матрица, что и в (16.69), а $\mathbf{F}_N = (\mathbf{F}_{1/2}, \dots, \mathbf{F}_{N-1/2})^T$, где

$$\mathbf{F}_{i-1/2} = \frac{1}{h_i} \int_{x_{i-1}}^{x_i} \left[\int_a^b K(x, s) f(s) ds \right] dx, \quad 1 \leq i \leq N.$$

Как мы уже видели, если выполнено условие (16.24), то $\|\mathbf{B}_N\|_{\infty} \leq q < 1$, и поэтому система уравнений (16.71) однозначно разрешима, а ее решение может быть найдено методом простой итерации

$$\mathbf{z}_N^{(k+1)} = \mathbf{B}_N \mathbf{z}_N^{(k)} + \mathbf{F}_N, \quad k \geq 0$$

при произвольном начальном приближении $\mathbf{z}_N^{(0)}$.

Если приближение (16.70) по методу Канторовича найдено, можно воспользоваться итерированным методом Канторовича (16.67) и найти уточненное приближение

$$\tilde{u}^h(x) = \sum_{j=1}^N \int_{x_{j-1}}^{x_j} K(x, s) ds \cdot z_{j-1/2} + \int_a^b K(x, s) f(s) ds + f(x).$$

5. Обоснование проекционных методов в случае использования кусочно-постоянных базисных функций. Приведем без доказательства две теоремы, дающие обоснование рассмотренных в п. 4 проекционных методов. Положим

$$L = \max_{a \leq x \leq b} \int_a^b |K'_x(x, s)| ds, \quad M = \max_{a \leq x \leq b} \int_a^b |K'_s(x, s)| ds.$$

Теорема 16.11. Пусть ядро K непрерывно на Π и выполнено условие (16.24), а функция f интегрируема по Риману на $[a, b]$. Тогда система уравнений метода Галеркина (16.69) имеет и притом единственное решение. Кроме того, метод Галеркина сходится:

$$\|u - u^h\| \rightarrow 0, \quad \text{при } h_{\max} \rightarrow 0.$$

Если же ядро K непрерывно дифференцируемо на Π , а функция f непрерывно дифференцируема на $[a, b]$, то для метода Галеркина верна оценка погрешности

$$\sup_{a \leq x \leq b} |u(x) - u^h(x)| \leq \frac{1}{1-q} \left(\frac{L}{1-q} \|f\|_{C[a, b]} + \|f'\|_{C[a, b]} \right) h_{\max},$$

а для итерированного метода Галеркина верна оценка погрешности

$$\|u - \tilde{u}^h\|_{C[a, b]} \leq \frac{M}{1-q} \left(\frac{L}{1-q} \|f\|_{C[a, b]} + \|f'\|_{C[a, b]} \right) h_{\max}^2. \blacksquare$$

Теорема 16.12. Пусть ядро K непрерывно на Π и выполнено условие (16.24), а функция f интегрируема по Риману на $[a, b]$. Тогда система уравнений метода Канторовича (16.70) имеет и притом единственное решение. Кроме того, метод Канторовича сходится:

$$\|u - u^h\|_{C[a, b]} \rightarrow 0 \quad \text{при } h_{\max} \rightarrow 0.$$

Если же ядро K непрерывно дифференцируемо на Π , то для метода Канторовича верна оценка погрешности

$$\|u - u^h\|_{C[a, b]} \leq \frac{L}{(1-q)^2} \sup_{a \leq x \leq b} |f(x)| h_{\max},$$

для интерированного метода Канторовича верна оценка погрешности

$$\|\tilde{u} - \tilde{u}^h\|_{C[a, b]} \leq \frac{ML}{(1-q)^2} \sup_{a \leq x \leq b} |f(x)| h_{\max}^2. \blacksquare$$

§ 16.6. Интегральные уравнения Фредгольма второго рода. Методы наименьших квадратов и коллокации

Рассмотрим кратко метод наименьших квадратов и метод коллокации, которые также используются для численного решения уравнений Фредгольма второго рода

1. Метод наименьших квадратов. Из неравенств (16.28) и (16.29) следует, что чем меньше будет невязка приближенного решения, тем больше оснований ожидать, что малой будет и величина погрешности. Напомним еще, что равенство нулю невязки означает, что найденное приближение является точным решением интегрального уравнения.

В *методе наименьших квадратов* предлагается искать приближенное решение в виде (16.58) или в виде (16.63) и определять вектор коэффициентов $\alpha_N = (\alpha_1, \alpha_2, \dots, \alpha_N)^T$ из условия минимума квадрата нормы невязки $\|r[u^N]\|$.

В первом случае невязка имеет вид (см. (16.59)):

$$r[u^N] = \sum_{j=1}^N \alpha_j \psi_j - f,$$

где $\psi_j = \varphi_j - \mathcal{K}\varphi_j$, и минимизируется величина

$$\|r[u^N]\|^2 = \left\| \sum_{j=1}^N \alpha_j \psi_j - f \right\|^2. \quad (16.72)$$

Во втором случае невязка имеет вид (см. (16.64)):

$$r[u^N] = \sum_{j=1}^N \alpha_j \psi_j - \mathcal{K}f,$$

и минимизируется величина

$$\|r[u^N]\|^2 = \left\| \sum_{j=1}^N \alpha_j \psi_j - \mathcal{K}f \right\|^2. \quad (16.73)$$

Заметим, что величина (16.72) представляет собой квадратичную функцию переменных $\alpha_1, \alpha_2, \dots, \alpha_N$:

$$\Phi(\alpha_1, \alpha_2, \dots, \alpha_N) = \sum_{k=1}^N \sum_{j=1}^N (\psi_k, \psi_j) \alpha_k \alpha_j - 2 \sum_{j=1}^N (f, \psi_j) \alpha_j + \|f\|^2.$$

Необходимое условие экстремума

$$\frac{\partial}{\partial \alpha_i} \Phi(\alpha_1, \alpha_2, \dots, \alpha_N) = 0, \quad 1 \leq i \leq N$$

приводит к системе линейных алгебраических уравнений

$$A_N \alpha_N = b_N, \quad (16.74)$$

где A_N — квадратная матрица с элементами $a_{ij} = (\psi_i, \psi_j)$, $1 \leq i, j \leq N$, а b_N — вектор-столбец с элементами $b_i = (f, \psi_i)$, $1 \leq i \leq N$.

Замечание. Обратим внимание на то, что матрица A — симметричная.

Аналогичным образом задача минимизации величины (16.73) сводится к системе вида (16.74) с той же матрицей A_N , но с другим вектором b_N , элементы которого вычисляются по формулам $b_i = (\mathcal{K}f, \psi_i)$, $1 \leq i \leq N$.

Для вычисления коэффициентов $\alpha_1, \alpha_2, \dots, \alpha_N$ можно воспользоваться одним из методов решения системы (16.74) либо одним из численных методов минимизации соответствующей квадратичной функции.

2. Метод коллокации. Согласно *методу коллокации* приближенное решение ищется в виде (16.58) или в виде (16.63). Коэффициенты линейной комбинации определяются так, чтобы невязка $r[u^N]$ обратилась в нуль в N выбранных на отрезке $[a, b]$ точках $a \leq x_1 < x_2 < \dots < x_N \leq b$, которые называются *точками коллокации*.

Поскольку для приближения (16.58) невязка имеет вид

$$r[u^N](x) = \sum_{j=1}^N \varphi_j(x) \alpha_j - \sum_{j=1}^N \int_a^b K(x, s) \varphi_j(s) ds \cdot \alpha_j - f(x),$$

то метод коллокации приводит к системе линейных алгебраических уравнений

$$\sum_{j=1}^N \Phi_j(x_i) \alpha_j = \sum_{j=1}^N \int_a^b K(x_i, s) \phi_j(s) ds \cdot \alpha_j + f(x_i), \quad i = 1, 2, \dots, N.$$

Если же используется приближение (16.63), то

$$r[u^N](x) = \sum_{j=1}^N \Phi_j(x) \alpha_j - \sum_{j=1}^N \int_a^b K(x, s) \phi_j(s) ds \cdot \alpha_j - \int_a^b K(x, s) f(s) ds,$$

и метод приводит к системе

$$\sum_{j=1}^N \Phi_j(x_i) \alpha_j = \sum_{j=1}^N \int_a^b K(x_i, s) \phi_j(s) ds \cdot \alpha_j + \int_a^b K(x_i, s) f(s) ds, \quad i = 1, 2, \dots, N.$$

§ 16.7. Интегральные уравнения Фредгольма второго рода. Метод замены ядра на вырожденное

1. Интегральные уравнения с вырожденным ядром. Ядро интегрального уравнения

$$\tilde{u}(x) = \int_a^b \tilde{K}(x, s) \tilde{u}(s) ds + f(x), \quad x \in [a, b] \quad (16.75)$$

называется *вырожденным*, если оно может быть представлено в виде конечной суммы вида

$$\tilde{K}(x, s) = \sum_{i=1}^N \Phi_i(x) \Psi_i(s). \quad (16.76)$$

Предполагается, что функции $\Phi_1, \Phi_2, \dots, \Phi_N$ линейно независимы на $[a, b]$; функции $\Psi_1, \Psi_2, \dots, \Psi_N$ также линейно независимы на $[a, b]$.

Уравнения с вырожденными ядрами решаются достаточно просто. После подстановки формулы (16.76) в уравнение (16.75) оно принимает вид

$$\tilde{u}(x) = \sum_{i=1}^N \Phi_i(x) \int_a^b \Psi_i(s) \tilde{u}(s) ds + f(x), \quad x \in [a, b]. \quad (16.77)$$

Следовательно, его решение имеет вид

$$\tilde{u}(x) = \sum_{i=1}^N c_i \Phi_i(x) + f(x), \quad (16.78)$$

где $c_i = \int_a^b \psi_i(s) \tilde{u}(s) ds$, $1 \leq i \leq N$

Подставляя выражение (16.78) в (16.77), получаем

$$\sum_{i=1}^N c_i \phi_i(x) = \sum_{i=1}^N \phi_i(x) \int_a^b \psi_i(s) \left[\sum_{j=1}^N c_j \phi_j(s) + f(s) \right] ds, \quad x \in [a, b].$$

Вследствие линейной независимости функций $\phi_1, \phi_2, \dots, \phi_N$ эти равенства верны тогда и только тогда, когда

$$c_i = \sum_{j=1}^N \int_a^b \psi_i(s) \phi_j(s) ds \cdot c_j + \int_a^b \psi_i(s) f(s) ds, \quad i = 1, 2, \dots, N.$$

Таким образом, для вычисления вектора неизвестных $c_N = (c_1, c_2, \dots, c_m)^T$ следует решить систему линейных алгебраических уравнений

$$c_N = B_N c_N + d_N,$$

где B_N — матрица с элементами $b_{ij} = \int_a^b \psi_i(s) \phi_j(s) ds$, $1 \leq i, j \leq N$, а d_N — вектор-столбец с элементами $d_i = \int_a^b \psi_i(s) f(s) ds$, $1 \leq i \leq N$.

После того как коэффициенты c_1, c_2, \dots, c_m вычислены, решение \tilde{u} может быть найдено по формуле (16.78).

2. Метод замены ядра на вырожденное. Пусть ядро K интегрального уравнения

$$u(x) = \int_a^b K(x, s) u(s) ds + f(x), \quad x \in [a, b] \quad (16.79)$$

может быть заменено на вырожденное ядро \tilde{K} вида (16.76), причем погрешность $K - \tilde{K}$ приближения достаточно мала. Тогда изложенным в п. 1 методом можно найти решение \tilde{u} интегрального уравнения (16.75) с вырожденным ядром \tilde{K} и рассматривать его как приближение к решению уравнения (16.79).

Как следует из неравенств (16.32) и (16.33), если погрешность $K - \tilde{K}$ замены ядра на вырожденное достаточно мала, то погрешность $u - \tilde{u}$ также будет мала.

Конечно, метод замены ядра на вырожденное может быть применен только тогда, когда приближение ядра вырожденным не является сложной задачей. Если ядро достаточно гладкое, а отрезок $[a, b]$ достаточно мал, то можно использовать конечный отрезок ряда Тейлора:

$$K(x, s) \approx \tilde{K}(x, s) = \sum_{n=0}^N \frac{(x - x_0)^n}{n!} \frac{\partial^n}{\partial x^n} K(x_0, s),$$

где $x \in [a, b]$.

Для построения вырожденного ядра можно попробовать использовать конечный отрезок двойного ряда Тейлора

$$K(x, s) \approx \tilde{K}(x, s) = \sum_{n=0}^N \sum_{m=0}^N \frac{1}{n! m!} \frac{\partial^{n+m}}{\partial x^n \partial s^m} K(x_0, s_0) (x - x_0)^n (s - s_0)^m$$

где $x_0, s_0 \in [a, b]$.

Иногда используют конечный отрезок разложения ядра по переменной x в ряд Фурье по тригонометрической системе функций или конечный отрезок двойного ряда Фурье. Можно заменять ядро вырожденным, используя методы интерполяции. Существуют и другие подходы.

Использование метода, основанного на приближении ядра вырожденным может оказаться значительно более трудоемким, чем численное решение интегрального уравнения методом квадратур или проекционными методами. В большинстве реальных задач, скорее всего, так оно и будет. Как правило, этот метод применяют тогда, когда решение нужно получить с невысокой точностью.

§ 16.8. Интегральные уравнения Вольтерра второго рода

1. Введение. Рассмотрим интегральное уравнение Вольтерра второго рода

$$u(x) = \int_a^x K(x, s) u(s) ds + f(x), \quad x \in [a, b] \quad (16.80)$$

Будем предполагать, что ядро K является непрерывной функцией, заданной на треугольнике $\Delta = \{(x, s) \mid a \leq x \leq b, a \leq s \leq x\}$, а функция f задана и непрерывна на отрезке $[a, b]$. Под решением интегрального уравнения (16.80) будем понимать непрерывную на отрезке $[a, b]$ функцию, удовлетворяющую уравнению для всех $x \in [a, b]$.

Формально уравнения Вольтерра можно рассматривать как частный случай уравнений Фредгольма. Достаточно доопределить ядро значением $K(x, s) = 0$ при $s > x$, чтобы можно было записать уравнение (16.80) в виде (16.79)¹. Но в действительности уравнения Фредгольма и Вольтерра обладают принципиально разными свойствами. Уравнения типа Вольтерра обладают свойством эволюционности: значение решения u в точке x полностью определяется его значениями в предшествующих точках (т.е. значениями в точках из промежутка $[a, x]$ — «предысторией» процесса); значения же в последующих точках (то есть в точках из промежутка $(x, b]$) никак не влияют на значение решения в точке x . В тех прикладных задачах, где появляются уравнения типа Вольтерра, переменная x , как правило, имеет смысл времени.

Напомним (см. § 16.2), что функция $y(t)$ является решением задачи Коши

$$y'(t) = f(t, y(t)),$$

$$y(t_0) = y_0$$

тогда и только тогда, когда она является решением нелинейного интегрального уравнения Вольтерра второго рода

$$y(t) = \int_0^t f(s, y(s)) \, ds + y_0.$$

Неудивительно, что теория уравнений Вольтерра второго рода во многом схожа с теорией задачи Коши для дифференциальных уравнений первого порядка. Есть значительное сходство и в численных методах решения.

2. Введение в теорию. Справедлива следующая теорема.

Теорема 16.13. Пусть ядро K непрерывно на множестве Δ , а функция f непрерывна на отрезке $[a, b]$. Тогда решение уравнения (16.80) существует, единственно и удовлетворяет оценке

$$\|u\|_{C[a, b]} \leq C_0 \|f\|_{C[a, b]} \quad (16.81)$$

с постоянной $C_0 = e^{M_0}$, где $M_0 = \int_a^b \max_{x \in [s, b]} |K(x, s)| \, ds$. ■

¹ Обратим внимание на то, что так доопределенное ядро перестает быть непрерывным.

Таким образом, задача решения интегрального уравнения Вольтерра второго рода корректна. Постоянная C_0 , входящая в оценку (16.81), играет в этой задаче роль числа обусловленности.

Гладкость решения уравнения (16.81) определяется гладкостью ядра K по переменной x , а также гладкостью правой части f . Важную роль играют также свойства «диагональной части» ядра — функции

$$\hat{K}_0(x) = K(x, x) \text{ и функция } \hat{K}_m(x) = \frac{\partial^m}{\partial x^m} K(x, s)|_{s=x}.$$

Теорема 16.14. *Пусть ядро K непрерывно на множестве Δ и имеет непрерывные на Δ частные производные $\frac{\partial K}{\partial x}$, $\frac{\partial^2 K}{\partial x^2}$, ..., $\frac{\partial^n K}{\partial x^n}$. Пусть функции \hat{K}_m ($0 \leq m \leq n$) являются $n - m - 1$ раз непрерывно дифференцируемыми на отрезке $[a, b]$, а правая часть f — n раз непрерывно дифференцируема на $[a, b]$. Тогда решение и уравнения (16.80) является n раз непрерывно дифференцируемой на отрезке $[a, b]$ функцией. ■*

3. Метод квадратур. Для решения уравнения (16.80) можно использовать метод квадратур. Выберем на отрезке $[a, b]$ узлы $a = x_0 < x_1 < \dots < x_N = b$. Введем шаги $h_n = x_n - x_{n-1}$, $1 \leq n \leq N$ и положим $h_{\max} = \max_{1 \leq n \leq N} h_n$.

Подставляя в уравнение (16.80) значения $x = x_n$ для $n = 0, 1, \dots, N$, приходим к системе равенств

$$u(x_0) = f(x_0),$$

$$u(x_n) = \int_a^{x_n} K(x_n, s) u(s) ds + f(x_n), \quad n = 1, 2, \dots, N.$$

Возьмем за основу квадратурные формулы вида

$$\int_a^{x_n} g(s) ds \approx \sum_{j=0}^n c_{nj} g(x_j),$$

с узлами x_0, x_1, \dots, x_n и весами $c_{nj} > 0$, которые могут зависеть от n .

Используя эти квадратурные формулы для аппроксимации интегралов $\int_a^{x_n} K(x_n, s) u(s) ds$, приходим к системе приближенных равенств

$$u(x_0) = f(x_0),$$

$$u(x_n) \approx \sum_{j=0}^n c_{nj} K_{nj} u(x_j) + f_n, \quad n = 1, 2, \dots, N.$$

Здесь $K_{nj} = K(x_n, x_j)$, $0 \leq j \leq n \leq N$ и $f_n = f(x_n)$, $0 \leq n \leq N$.

Потребуем теперь, чтобы приближения $u_n = u^N(x_n)$ к значениям $u(x_n)$ в точках x_n удовлетворяли системе линейных алгебраических уравнений

$$u_0 = f_0,$$

$$u_n = \sum_{j=0}^n c_{nj} K_{nj} u_j + f_n, \quad n = 1, 2, \dots, N.$$

Если $c_{nn} K_{nn} \neq 1$ для всех n , то значения неизвестных u_0, u_1, \dots, u_N могут быть последовательно найдены по формулам

$$u_0 = f_0,$$

$$u_n = \frac{1}{1 - c_{nn} K_{nn}} \left[\sum_{j=0}^{n-1} c_{nj} K_{nj} u_j + f_n \right], \quad n = 1, 2, \dots, N.$$

Приближенное решение u^N можно вычислить и в произвольной промежуточной точке из отрезка $[a, b]$. Если $x \in (x_{n-1}, x_n)$, то для аппроксимации интеграла $\int_a^x K(x, s) u(s) ds$ в этом случае следует использовать квадратурную формулу

$$\int_a^x g(s) ds \approx \sum_{j=0}^{n-1} c_{xj} g(x_j) + c_{xn} g(x)$$

с узлами $x_0, x_1, \dots, x_{n-1}, x$. Весы $c_{x0}, c_{x1}, \dots, c_{xn}$ при этом, естественно, могут зависеть от x . Формула для вычисления $u^N(x)$ примет в этом случае вид

$$u^N(x) = \frac{1}{1 - c_{xn} K(x, x)} \left[\sum_{j=0}^{n-1} c_{xj} K(x, x_j) u_j + f(x) \right], \quad x \in (x_{n-1}, x_n).$$

Конечно, для вычисления решения в промежуточных точках, можно воспользоваться и одним из методов интерполяции.

4. Применение формулы трапеций. Покажем, как для вычисления приближенного решения u^N уравнения (16.80) может быть использована составная формула трапеций

$$\int_a^{x_n} g(s) \, ds \approx g(x_0) \frac{h_1}{2} + \sum_{j=1}^{n-1} g(x_j) h_{j+1/2} + g(x_n) \frac{h_n}{2},$$

где $h_{j+1/2} = (h_j + h_{j+1})/2$.

Используя ее для аппроксимации интеграла $\int_a^{x_n} K(x_n, s) u(s) \, ds$, приходим к системе уравнений

$$u_0 = f_0,$$

$$u_1 = \frac{K_{10}h_1}{2} u_0 + \frac{K_{11}h_1}{2} u_1 + f_1,$$

$$u_n = \frac{K_{n0}h_1}{2} u_0 + \sum_{j=1}^{n-1} K_{nj} h_{j+1/2} u_j + \frac{K_{nn}h_n}{2} u_n + f_n, \quad n = 2, 3, \dots, N$$

относительно значений $u_n = u^N(x_n) \approx u(x_n)$, $0 \leq n \leq N$.

Если $K_{nn}h_n \neq 2$ для всех n , то значения u_0, u_1, \dots, u_N могут быть последовательно найдены по формулам

$$u_0 = f_0,$$

$$u_1 = \frac{\frac{K_{10}h_1}{2} u_0 + f_1}{1 - K_{11} \frac{h_1}{2}}, \quad (16.82)$$

$$u_n = \frac{\frac{K_{n0}h_1}{2} u_0 + \sum_{j=1}^{n-1} K_{nj} h_{j+1/2} u_j + f_n}{1 - K_{nn} \frac{h_n}{2}}, \quad n = 2, 3, \dots, N. \quad (16.83)$$

Для того чтобы найти значения приближения $u^N(x_n)$ в промежуточных точках $x \in (x_{n-1}, x_n)$, используем составную квадратурную формулу трапеций с узлами $x_0, x_1, \dots, x_{n-1}, x$:

$$\int_a^x g(s) \, ds \approx g(x_0) \frac{h_1(x)}{2} + g(x) \frac{h_1(x)}{2} \quad \text{при } n = 1,$$

$$\int_a^x g(s) \, ds \approx g(x_0) \frac{h_1}{2} + \sum_{j=1}^{n-2} g(x_j) h_{j+1/2} + g(x_{n-1}) h_{n-1/2}(x) + g(x) \frac{h_n(x)}{2}$$

при $2 \leq n \leq k$. Здесь $h_n(x) = x - x_{n-1}$, $h_{n-1/2}(x) = (h_{n-1} + h_n(x))/2$.

В результате мы придем к формулам

$$u^N(x) = \frac{K(x, x_0) \frac{h_1(x)}{2} u_0 + f(x)}{1 - K(x, x) \frac{h_1(x)}{2}}, \quad \text{для } x \in (x_0, x_1), \quad (16.84)$$

$$u^N(x) = \frac{K(x, x_0) \frac{h_1}{2} u_0 + \sum_{j=1}^{n-2} K(x, x_j) h_{j+1/2} u_j + K(x, x_{n-1}) h_{n-1/2}(x) u_{n-1}}{1 - K(x, x) \frac{h_n(x)}{2}} + \\ + \frac{f(x)}{1 - K(x, x) \frac{h_n(x)}{2}} \quad \text{для } x \in (x_{n-1}, x_n), \quad 2 \leq n \leq N. \quad (16.85)$$

Заметим, что в случае, когда $\hat{K}_{0,\max} = \max_{x \in [a, b]} \hat{K}_0(x) \leq 0$ (т.е. когда

$\hat{K}_{0,\max} = \hat{K}(x, x) \leq 0$ для всех $x \in [a, b]$), знаменатели в расчетных формулах (16.82)–(16.85) не обращаются в ноль и все вычисления можно провести до конца. Однако в случае, когда $\hat{K}_{0,\max} > 0$, для этого следует потребовать достаточную малость максимального из шагов:

$$\hat{K}_{0,\max} h_{\max} \leq 1. \quad (16.86)$$

Приведем без доказательства теорему об оценке погрешности рассмотренного в этом пункте метода.

Теорема 16.15. Пусть ядро K является непрерывной на множестве Δ функцией и имеет непрерывные на Δ производные K'_x , K''_{xx} . Пусть функция \hat{K}_0 непрерывно дифференцируема на отрезке $[a, b]$, а функция f дважды непрерывно дифференцируема на $[a, b]$. Пусть также выполнено условие (16.86). Тогда справедлива оценка погрешности

$$\|u - u^N\|_{C[a, b]} \leq Ch_{\max}^2,$$

$$\text{где } C = e^{2(b-a)M} \frac{b-a}{12} \max_{(x, s) \in \Delta} \left| \frac{\partial^2}{\partial x^2} [K(x, s)u(s)] \right|, \quad M = \max_{(x, s) \in \Delta} |K(x, s)|. \blacksquare$$

§ 16.9. Интегральные уравнения Фредгольма первого рода

Задача решения интегрального уравнения Фредгольма первого рода

$$\int_a^b K(x, s) u(s) \, ds = f(x), \quad x \in [a, b] \quad (16.87)$$

относится к классу некорректных задач. Это проявляется в том, что уравнение (16.87) не может быть разрешено при произвольной правой части f . Если же при специальным образом выбранной правой части f уравнение имеет решение, то это решение может и не быть единственным. Более того, сколь угодно малые погрешности задания правой части f могут приводить к сколь угодно большим погрешностям решения.

Приведем простой пример, иллюстрирующий проблемы, возникающие при решении уравнения Фредгольма первого рода. Ситуация в общем случае более сложная, но сходная.

Пример 16.1. Рассмотрим уравнение (16.87) с ядром $K(x, s) \equiv 1$, то есть уравнение

$$\int_a^b u(s) \, ds = f(x), \quad x \in [a, b]. \quad (16.88)$$

Заметим, что интегральный оператор, стоящий в левой части уравнения (16.88), преобразует всякую функцию u в константу. Это означает, что уравнение заведомо не имеет решений, если $f(x) \not\equiv \text{const}$.

Если же $f(x) \equiv c = \text{const}$, то решением уравнения является $u(x) = \frac{c}{b-a}$. Но решением того же уравнения является и любая функция вида

$$u(x) = \frac{c}{b-a} + \varphi(x), \quad \text{где } \int_a^b \varphi(x) \, dx = 0.$$

Таким образом, множество решений уравнения (16.88) бесконечно, и в нем есть решения, сколь угодно сильно отличающиеся друг от друга.

Чтобы сделать задачу решения уравнения (2) корректной, нужно ограничить множество правых частей уравнения постоянными и искать решение также в классе постоянных.▲

1. Метод регуляризации Тихонова. Для решения некорректных задач разработаны различные подходы. Мы кратко рассмотрим *метод регуляризации*, разработанный А.Н. Тихоновым, на примере решения интегрального уравнения Фредгольма первого рода.

Пусть требуется найти решение уравнения (16.87), т.е. уравнения

$$\mathcal{K}u = f. \quad (16.89)$$

При этом вместо идеальной правой части f задано ее приближение \tilde{f} , так что фактически приходится решать уравнение

$$\mathcal{K}u = \tilde{f}. \quad (16.90)$$

Предположим для простоты, что интегральный оператор устроен так, что однородное уравнение $\mathcal{K}u = 0$ имеет только нулевое решение. В этом случае неоднородное уравнение (16.89) при заданной части f имеет решение и оно единственno. Обозначим точное решение уравнения (16.89), отвечающее правой части f , через $u[f]$.

Рассматриваемая задача некорректна. Поэтому, как бы не была мала погрешность $\tilde{f} - f$ задания правой части, погрешность решения $u[f] - u[\tilde{f}]$ может оказаться сколь угодно большой.

Таким образом, пытаться искать точное решение интегрального уравнения (16.90) бессмысленно. В методе регуляризации предлагается заменить исходную задачу другой, *регуляризованной задачей*, решение которой близко к решению исходной задачи. Как правило, регуляризованныя задача является корректной, но плохо обусловленной. Один из способов построения такой регуляризованной задачи излагается ниже.

Обычно известен уровень погрешности задания правой части. Будем предполагать, что

$$\|f - \tilde{f}\| < \delta,$$

где $0 < \delta$ — величина, характеризующая уровень погрешности.

Определим невязку, отвечающую приближенному решению \tilde{u} уравнения (16.89) формулой $r[\tilde{u}] = \mathcal{K}\tilde{u} - f$. Ясно, что функция u является точным решением интегрального уравнения (16.89) тогда и только тогда, когда $r[u] = 0$. Поэтому решение уравнения является точкой минимума функционала

$$J(u) = \|r[u]\|^2 = \|\mathcal{K}u - f\|^2$$

Так что задача минимизации функционала $J(u)$ эквивалентна задаче: решения интегрального уравнения (16.89).

Положим $\tilde{J}(u) = \|\mathcal{K}u - \tilde{f}\|^2$ и введем в функционал дополнительное слагаемое $\Omega_\alpha = \alpha\|u\|^2$ (это слагаемое принято называть *регуляризатором*), где $0 < \alpha$ — параметр (*параметр регуляризации*). Будем теперь рассматривать регуляризованную задачу — задачу минимизации функционала

$$J_\alpha(u) = \tilde{J}(u) + \Omega_\alpha(u) = \|\mathcal{K}u - \tilde{f}\|^2 + \alpha\|u\|^2.$$

Известно, что при любой правой части \tilde{f} регуляризованная задача имеет единственное решение, которое мы обозначим через $u_\alpha[\tilde{f}]$.

Можно также показать, что регуляризованная задача эквивалентна задаче решения интегрального уравнения Фредгольма второго рода

$$u(x) = \int_a^b K_\alpha(x, s)u(s) ds + F_\alpha(x),$$

где

$$K_\alpha(x, s) = -\alpha^{-1} \int_a^b K(x, t)K(t, s) dt, \quad F_\alpha(x) = \alpha^{-1} \int_a^b K(x, s)\tilde{f}(s) ds.$$

Более того, эта задача устойчива, а именно, справедлива оценка

$$\|u_\alpha[f] - u_\alpha[\tilde{f}]\| \leq C(\alpha)\|f - \tilde{f}\|,$$

в которой положительная постоянная $C(\alpha)$ не зависит от f и \tilde{f} .

Если величина параметра регуляризации $\alpha = \alpha(\delta)$ правильно согласована с величиной погрешности δ , то за приближенное решение уравнения (16.89) можно принять решение $u_{\alpha(\delta)}[\tilde{f}]$ регуляризованной задачи, причем приближенное решение будет сходиться к точному при $\delta \rightarrow 0$.

Теорема 16.16. Пусть $\alpha = \alpha(\delta) > 0$, причем $\alpha(\delta) \rightarrow 0$ и $\frac{\delta^2}{\alpha(\delta)} \rightarrow 0$

при $\delta \rightarrow 0$. Тогда $\|u_{\alpha(\delta)}[\tilde{f}] - u[f]\| \rightarrow 0$ ■

Замечание. Метод регуляризации применим и в том случае, когда решение уравнения (16.89) не является единственным. В этом случае сходимость будет иметь место к так называемому *нормальному решению*, т.е. к решению уравнения (16.89) с минимальной нормой $\|u\|$.

Для выбора конкретной величины параметра регуляризации $\alpha(\delta)$, часто используется метод выбора параметра регуляризации по невязке. Параметр α определяется как решение нелинейного уравнения

$$\varphi(\alpha) = \delta^2,$$

где $\varphi(\alpha) = \| \mathcal{K}u_\alpha[\tilde{f}] - \tilde{f} \|^2$.

Часто известно, что решение уравнения (16.89) — гладкое. В этой ситуации используется регуляризатор

$$\Omega_\alpha = \alpha [\|u\|^2 + \|u'\|^2],$$

и регуляризованной задачей становится задача поиска точки минимума функционала

$$J_\alpha(u) = \|\mathcal{K}u - \tilde{f}\|^2 + \alpha [\|u\|^2 + \|u'\|^2]. \quad (16.91)$$

В этом случае, если $\alpha(\delta) \rightarrow 0$, $\frac{\delta^2}{\alpha(\delta)} \rightarrow 0$ при $\delta \rightarrow 0$, то будет выполнено не только свойство $\|u_{\alpha(\delta)}[\tilde{f}] - u[f]\| \rightarrow 0$, но и дополнительное свойство $\|u'_{\alpha(\delta)}[\tilde{f}] - u'[f]\| \rightarrow 0$.

2. Дискретизация. Для того чтобы численно реализовать метод регуляризации, необходимо перейти к его дискретному варианту.

Покажем, как это можно сделать, используя метод квадратур. Пусть для аппроксимации интегрального оператора используется формула прямоугольников (16.45) с достаточно малыми шагами h_j , а производная u' аппроксимируется центральной разностной производной

$$u'(x_i) \approx \frac{u_{i+1/2} - u_{i-1/2}}{h_{i+1/2}}, \quad 0 < i < N.$$

Тогда дискретным вариантом функционала $\tilde{J}(u)$ будет функционал

$$\tilde{J}^h(u^h) = \sum_{i=1}^N \left[\sum_{j=1}^N K_{i-1/2, j-1/2} h_j u_{j-1/2} - \tilde{f}(x_{i-1/2}) \right]^2 h_i,$$

дискретным аналогом функционала $\Omega_\alpha(u) = \alpha [\|u\|^2 + \|u'\|^2]$ — функционал

$$\Omega_\alpha^h(u^h) = \alpha \left[\sum_{i=1}^N u_{i-1/2}^2 h_i + \sum_{i=1}^{N-1} \left(\frac{u_{i+1/2} - u_{i-1/2}}{h_{i+1/2}} \right)^2 h_{i+1/2} \right],$$

а дискретным аналогом функционала (16.91) — функционал

$$\begin{aligned} J_\alpha^h(u^h) = \tilde{J}^h(u^h) + \Omega_\alpha^h(u^h) &= \sum_{i=1}^N \left[\sum_{j=1}^N K_{i-1/2, j-1/2} h_j u_{j-1/2} - \tilde{f}(x_{i-1/2}) \right]^2 h_i + \\ &+ \alpha \left[\sum_{i=1}^N u_{i-1/2}^2 h_i + \sum_{i=1}^{N-1} \left(\frac{u_{i+1/2} - u_{i-1/2}}{h_{i+1/2}} \right)^2 h_{i+1/2} \right]. \end{aligned}$$

Отметим, что задача минимизации этого функционала эквивалентна линейной задаче метода наименьших квадратов (см. § 5.12), примененной к решению переопределенной системы линейных алгебраических уравнений

$$\sqrt{h_i} \sum_{j=1}^N K_{i-1/2, j-1/2} h_j u_{j-1/2} = \sqrt{h_i} \tilde{f}(x_{i-1/2}), \quad 1 \leq i \leq N,$$

$$\sqrt{\alpha} \sqrt{h_i} u_{i-1/2} = 0, \quad 1 \leq i \leq N,$$

$$\sqrt{\alpha} \frac{u_{i+1/2} - u_{i-1/2}}{\sqrt{h_{i+1/2}}} = 0, \quad 1 \leq i < N.$$

§ 16.10. Интегральные уравнения Вольтерра первого рода

Рассмотрим интегральное уравнение Вольтерра первого рода

$$\int_a^x K(x, s) u(s) ds = f(x), \quad x \in [a, b]. \quad (16.92)$$

Заметим, что при $x = a$ функция, стоящая в левой части уравнения, обращается в нуль. Поэтому для разрешимости уравнения необходимо, чтобы выполнялось условие $f(a) = 0$.

Пусть, как и ранее, $\hat{K}_0(x) = K(x, x)$. Предположим, что ядро K имеет непрерывную частную производную K'_x . В этом случае левая часть уравнения является непрерывно дифференцируемой на отрезке $[a, b]$ функцией, причем

$$\frac{d}{dx} \int_a^x K(x, s) u(s) ds = \hat{K}_0'(x) u(x) + \int_a^x K'_x(x, s) u(s) ds.$$

Поэтому и правая часть f также должна быть непрерывно дифференцируемой на отрезке $[a, b]$.

Дифференцируя уравнение (16.92), получаем

$$K_0(x)u(x) + \int_a^x K'_x(s)u(s) ds = f'(x).$$

Если $\hat{K}_0(x)u(x) \neq 0$ на $[a, b]$, то полученное уравнение может быть записано в виде уравнения Вольтерра второго рода

$$u(x) = \int_a^x G(x, s)u(s) ds + F(x), \quad (16.93)$$

где $G(x, s) = -\frac{1}{\hat{K}_0(x)} K'_x(s)$, $F(x) = \frac{1}{\hat{K}_0(x)} f'(x)$.

Из результатов предыдущего параграфа вытекает следующий результат.

Теорема 16.17. Пусть ядро K и частная производная ядра K'_x непрерывны на множестве Δ . Предположим, что $\hat{K}_0(x) \neq 0$ на отрезке $[a, b]$.

Если функция f непрерывно дифференцируема на отрезке $[a, b]$ и удовлетворяет условию $f(a) = 0$, то решение интегрального уравнения (16.92) существует, единственно и удовлетворяет оценке

$$\max_{[a, b]} |u(x)| \leq C \max_{[a, b]} |f'(x)|,$$

где $C = \frac{1}{\mu} e^{M(b-a)}$, $\mu = \min_{[a, b]} |\hat{K}_0(x)|$, $M = \max_{(x, s) \in \Delta} \left| \frac{1}{\hat{K}_0(x)} K'_x(s) \right|$. ■

Таким образом, в условиях теоремы 16.16 задача решения интегрального уравнения Вольтерра первого рода с приближенно заданной правой частью \tilde{f} корректна, если для приближений \tilde{f} к f возможно гарантировать малость погрешности $(\tilde{f})' - f'$.

Решение уравнения (16.93) может быть найдено с помощью одного из численных методов решения уравнения Вольтерра второго рода, если значения производных K'_x и f' могут быть найдены с высокой точностью.

2. Некорректность задачи. Из изложенного в п. 1 материала могло сложиться впечатление, что численное решение интегральных уравнений Вольтерра первого рода — простая задача. Однако это не так.

Как правило, правая часть f уравнения задается приближенно, т.е. вместо уравнения (16.92) приходится решать уравнение

$$\int_a^x K(x,s)u(s) \, ds = \tilde{f}(x), \quad x \in [a, b], \quad (16.94)$$

где $\tilde{f}(x) \approx f(x)$. Если $\tilde{f}(a) \neq 0$ или функция \tilde{f} не является непрерывно дифференцируемой, то уравнение (16.94) не имеет решений. Кроме того, предложенный подход предполагает, что производная $(\tilde{f})' \approx f'$, что весьма маловероятно.

Часть возникающих проблем видна из следующего примера.

Пример 16.2. В случае, когда $K(x,s) \equiv 1$ уравнение (16.92) принимает вид

$$\int_a^x u(s) \, ds = f(x).$$

Ясно, что решением этого уравнения является функция u , являющаяся производной функции f :

$$u(x) = f'(x).$$

Таким образом, решение уравнения сводится к дифференцированию правой части уравнения. Как известно (см. § 3.1, пример 3.5), задача дифференцирования функции f по ее приближенно заданным значениям $\tilde{f}(x)$ является неустойчивой. ▲

Итак, задача вычисления решения интегрального уравнения Вольтерра первого рода с приближенно заданной правой частью некорректна. Ее грамотное численное решение предполагает использование специальных методов, например методов регуляризации.

Замечание. Обратим внимание на то, что предложенный в п. 1 подход неприменим для уравнений с негладкими ядрами. Например, его нельзя использовать для уравнения Абеля (см. § 16.2)

$$\int_a^x \frac{u(s) \, ds}{\sqrt{x-s}} = f(x),$$

поскольку его ядро $K(x,s) = \frac{1}{\sqrt{x-s}}$ не определено при $x = s$.

Кроме того, предположение о том, что $\hat{K}_0(x) \neq 0$, даже для гладких ядер выполняется далеко не всегда. Например, оно не выполнено для уравнения (16.15) с ядром $K(x,s) = \frac{(x-s)^{n-1}}{(n-1)!}$ при $n > 1$.

§ 16.11. Нелинейные интегральные уравнения

Рассмотренные в данной главе методы могут быть применены и к решению нелинейных интегральных уравнений Фредгольма и Вольтерра.

1. Нелинейное интегральное уравнение Фредгольма второго рода. Метод квадратур. Для численного решения уравнения

$$u(x) = \int_a^b K(x, s, u(s)) ds + f(x), \quad x \in [a, b] \quad (16.95)$$

может быть использован метод квадратур. Если за основу взята квадратурная формула (16.39), то метод квадратур приводит не к системе линейных алгебраических уравнений (16.41), как это было при решении линейного интегрального уравнения (16.38), а к системе нелинейных уравнений

$$u_i = \sum_{j=1}^N c_j K(x_i, x_j, u_j) + f_i, \quad i = 1, 2, \dots, N \quad (16.96)$$

относительно неизвестных u_1, u_2, \dots, u_N , где $u_i = u^N(x_i)$, $1 \leq i \leq N$.

Решив эту систему одним из итерационных методов, мы можем затем вычислить приближенное решение во всех точках отрезка $[a, b]$ по формуле

$$u^N(x) = \sum_{j=1}^N c_j K(x, x_j, u_j) + f(x), \quad x \in [a, b]. \quad (16.97)$$

Если, например, за основу взята квадратурная формула прямоугольников (16.45), то система (16.96) принимает вид

$$u_{i-1/2} = \sum_{j=1}^N K(x_{i-1/2}, x_{j-1/2}, u_{j-1/2}) h_j + f(x_{i-1/2}), \quad i = 1, 2, \dots, N,$$

где $u_{i-1/2} = u^N(x_{i-1/2})$, $1 \leq i \leq N$.

При этом формула (16.97) выглядит так:

$$u^N(x) = \sum_{j=1}^N K(x, x_{j-1/2}, u_{j-1/2}) h_j + f(x), \quad x \in [a, b].$$

2. Нелинейное интегральное уравнение Фредгольма второго рода. Метод Галеркина. Для решения уравнения (16.95) возможно и

применение проекционных методов. Классический метод Галеркина приводит к системе нелинейных уравнений

$$\sum_{j=1}^N \int_a^b \phi_j(x) \phi_i(x) dx \cdot a_j = \int_a^b \int_a^b K\left(x, s, \sum_{j=1}^N a_j \phi_j(x)\right) \phi_i(x) dx ds + \int_a^b f(x) \phi_i(x) dx, \quad 1 \leq i \leq N. \quad (16.98)$$

В частности, при использовании метода Галеркина с кусочно-постоянными базисными функциями система (16.98) принимает вид (ср. с (16.69)):

$$u_{i-1/2} = \Psi_i(u_{1/2}, \dots, u_{N-1/2}) + f_{i-1/2}, \quad 1 \leq i \leq N,$$

$$\text{где } \Psi_i(u_{1/2}, \dots, u_{N-1/2}) = \frac{1}{h_i} \int_{x_{i-1}}^{x_i} \left[\sum_{j=1}^N \int_{x_{j-1}}^{x_j} K(x, s, u_{j-1/2}) ds \right] dx, \\ f_{i-1/2} = \frac{1}{h_i} \int_{x_{i-1}}^{x_i} f(x) dx.$$

3. Нелинейное интегральное уравнение Вольтерра второго рода. Для численного решения уравнения

$$u(x) = \int_a^x K(x, s, u(s)) ds + f(x), \quad x \in [a, b] \quad (16.99)$$

также может быть использован метод квадратур. Если в случае линейного уравнения (16.80) задача сводилась к системе линейных алгебраических уравнений, то для нелинейного уравнения (16.99) метод квадратур дает систему нелинейных уравнений

$$u_0 = f_0,$$

$$u_n = \sum_{j=0}^n c_{nj} K(x_n, x_j, u_j) + f_n, \quad n = 1, 2, \dots, N.$$

Значения неизвестных u_0, u_1, \dots, u_N вычисляются последовательно. Если u_1, u_2, \dots, u_{n-1} уже найдены, то для вычисления очередного значения u_n нужно решить нелинейное уравнение

$$u_n = c_{nn} K(x_n, x_n, u_n) + F_n,$$

$$\text{где } F_n = \sum_{j=0}^{n-1} c_{nj} K(x_n, x_j, u_j) + f_n.$$

4. Нелинейное интегральное уравнение Фредгольма первого рода.

Задача вычисления решения уравнения

$$\int_a^b K(x, s, u(s)) ds = f(x), \quad x \in [a, b] \quad (16.99)$$

некорректна. Поэтому для ее решения необходимо использование методов регуляризации.

Например, по аналогии с изложенным в § 16.7 подходом эту задачу можно заменить регуляризованной задачей минимизации функционала

$$J_\alpha(u) = \int_a^b \left[\int_a^b K(x, s, u(s)) ds - \tilde{f}(x) \right]^2 dx + \alpha [\|u\|^2 + \|u'\|^2]$$

и для численной реализации использовать ее дискретный вариант — задачу минимизации функционала

$$\begin{aligned} J_\alpha^h(u^h) = & \sum_{i=1}^N \left[\sum_{j=1}^N K(x_{i-1/2}, x_{j-1/2}, u_{j-1/2}) h_j - \tilde{f}(x_{j-1/2}) \right]^2 h_i + \\ & + \alpha \left[\sum_{j=1}^N |u_{j-1/2}|^2 h_j + \sum_{j=1}^{N-1} \left(\frac{u_{j+1/2} - u_{j-1/2}}{h_{j+1/2}} \right)^2 h_{j+1/2} \right]. \end{aligned}$$

§ 16.12. Дополнительные замечания

1. Мы изложили только методы решения одномерных интегральных уравнений, в которых решением является функция одной переменной, а интегрирование производится по отрезку $[a, b]$. Однако на практике часто встречаются многомерные интегральные уравнения Фредгольма первого и второго рода

$$\int_G K(x, s) u(s) ds = f(x), \quad x \in G, \quad (16.100)$$

$$u(x) = \int_G K(x, s) u(s) ds + f(x), \quad x \in G, \quad (16.101)$$

в которых искомая функция $u(x) = u(x_1, x_2, \dots, x_m)$ является функцией многих переменных, а интегрирование проводится по области $G \subset \mathbb{R}^m$. Встречаются и задачи, в которых интегрирование производится по некоторым кривым или поверхностям.

Рассмотренные в данной главе методы могут быть применены (с естественными изменениями) и к численному решению уравнений (16.100), (16.101).

2. Для простоты изложения и понимания мы ограничились в данной главе интегральными уравнениями с вещественнонзначными ядрами K и правыми частями f . В приложениях нередко встречаются и интегральные уравнения с комплекснозначными ядрами K и правыми частями f . Теория и методы решения таких уравнений в основном аналогичны вещественному случаю. Значительная часть отличий при этом носит формальный характер. Так, скалярное произведение комплекснозначных функций f и g следует определять формулой

$$(f, g) = \int_a^b f(x) \overline{g(x)} dx.$$

(Как обычно, черта над комплексным числом $z = a + ib$ означает взятие операции комплексного сопряжения: $\bar{z} = a - ib$.)

Для уравнений с комплекснозначными ядрами роль сопряженного оператора \mathcal{K}^* играет оператор с ядром $K^*(x, s) = \overline{K(s, x)}$, а условие (16.17) симметричности ядра заменяется условием

$$K(x, s) = \overline{K(s, x)}.$$

3. Во многих прикладных задачах ядро K не является гладким. Оно может иметь те или иные особенности. Чаще всего ядро перестает быть гладким на диагонали $x = s$ квадрата Π . Например, ядро $G(x, s)$ уравнения (16.12) непрерывно, но его производные имеют разрыв при $x = s$.

Ядро $K(z, s) = \frac{1}{\sqrt{z-s}}$ уравнения Абеля (16.5) стремится к бесконечности при $z - s \rightarrow 0$. Ядро $E_1(|\tau - \tau'|)$ уравнения переноса излучения (16.6) также стремится к бесконечности при $\tau - \tau' \rightarrow 0$.

Для решения интегральных уравнений с такими ядрами проекционные методы и метод наименьших квадратов приспособлены значительно лучше, чем метод квадратур.

4. Одной из практически важных задач является задача вычисления собственных значений интегральных операторов.

Напомним, что комплексное число λ называется *собственным значением* интегрального оператора \mathcal{K} , если существует ненулевая функция u такая, что

$$\int_a^b K(x, s) u(s) \, ds = \lambda u(x), \quad x \in [a, b]. \quad (16.102)$$

Эта функция называется *собственной функцией* интегрального оператора, отвечающей собственному значению λ .

Известно, что всякий интегральный оператор \mathcal{K} с ядром, удовлетворяющим условию

$$\int_a^b \int_a^b |K(x, s)|^2 \, ds \, dx < \infty,$$

имеет конечный или счетный набор собственных значений. Если множество собственных значений счетно, то его можно упорядочить в порядке убывания модулей: $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n| \geq \dots$. При этом $\lambda_n \rightarrow 0$ при $n \rightarrow \infty$.

Известно также, что все собственные числа самосопряженного оператора \mathcal{K} являются вещественными, а его собственные функции, отвечающие различным собственным значениям, ортогональны.

Изложенные в данной главе методы могут быть использованы и для решения задачи вычисления собственных значений интегрального оператора \mathcal{K} . Например, аппроксимируя уравнение (16.102) в соответствии с изложенным в п. 1 § 16.3 методом квадратур, мы приходим к алгебраической задаче на собственные значения

$$\mathbf{B}_N \mathbf{u}_N = \lambda \mathbf{u}_N.$$

Собственные значения $\lambda_1^{(N)}, \lambda_2^{(N)}, \dots, \lambda_N^{(N)}$ матрицы \mathbf{B}_N дают приближения к первым собственным значениям $\lambda_1, \lambda_2, \dots, \lambda_N$ интегрального оператора \mathcal{K} .

Применяя к уравнению (16.102) метод Галеркина в изложенном в п. 1 § 16.4 варианте, мы приходим к обобщенной алгебраической задаче на собственные значения

$$\mathbf{K}_N \alpha = \lambda \mathbf{D}_N \alpha.$$

Собственные значения матрицы $\mathbf{D}_N^{-1} \mathbf{K}_N$ также дают приближения к собственным значениям оператора \mathcal{K} .

Следует иметь в виду, что лишь для первых M (где M значительно меньше N) наибольших по модулю собственных значений оператора \mathcal{K} получаются достаточно точные приближения.

5. Начала теории интегральных уравнений можно найти в учебниках [3*, 7*, 19*]. Аналитические методы решения ряда интегральных уравнений со специальными ядрами содержатся в справочнике [9*]. Разнообразная информация о численных методах решения интегральных уравнений содержится в учебниках [9, 14, 22, 36, 51, 61, 1*, 3*, 4*] и монографии [19*]. Более подробно методы решения некорректных интегральных уравнений изложены в [7, 37, 97, 98].

СПИСОК ЛИТЕРАТУРЫ

1. Алберг Дж., Нилсон Э., Уолш Дж. Теория сплайнов и ее приложения. М.: Мир, 1972.
2. Андреев В.Б., Руховец Л.А. Проекционные методы. М.: Знание, 1986.
3. Арушанян О.Б., Залеткин С.В. Численное решение обыкновенных дифференциальных уравнений на Фортране. М.: Изд-во МГУ, 1990.
4. Бабенко К.И. Основы численного анализа. М.: Наука, 1986.
5. Бабушка И., Витасек Э., Прагер М. Численные процессы решения дифференциальных уравнений. М.: Мир, 1969.
6. Базара М., Шетти К. Нелинейное программирование. Теория и алгоритмы. М.: Мир, 1982.
7. Бакушинский А.Б., Гончарский А.В. Некорректные задачи. Численные методы и приложения. М.: Изд-во МГУ, 1989.
8. Бахвалов Н.С. Численные методы. М.: Наука, 1973.
9. Бахвалов Н.С., Жидков Н.П., Кобельков Г.М. Численные методы. М.: Наука, 1987.
10. Бахвалов Н.С., Лапин А.В., Чижонков Е.В. Численные методы в задачах и упражнениях. М.: Высшая школа, 2000.
11. Бейкер Дж., Грейвс-Моррис П. Аппроксимации Паде. М.: Мир, 1986.
12. Беклемишев Д.В. Дополнительные главы линейной алгебры. М.: Наука, 1983.
13. Березин И.С., Жидков Н.П. Методы вычислений. Т. 1. М.: Физматгиз, 1962.
14. Березин И.С., Жидков Н.П. Методы вычислений. Т. 2. М.: Физматгиз, 1962.
15. Блехман Н.Н., Мышикис А.Д., Пановко Я.Г. Механика и прикладная математика. Логика и особенности приложений математики. М.: Наука, 1983.
16. Боглаев Ю.П. Вычислительная математика и программирование. М.: Высшая школа, 1990.
17. Боголюбов А.Н. Математики. Механики: Биографический справочник. Киев: Наукова думка, 1983.
18. Де Бор К. Практическое руководство по сплайнам. М.: Радио и связь, 1985.
19. Бородин А.Н., Бугай А.С. Выдающиеся математики. Киев: Радянська школа, 1987.
20. Васильев Ф.П. Численные методы решения экстремальных задач. М.: Наука, 1980.
21. Вержбицкий В.М. Численные методы. Линейная алгебра и нелинейные уравнения. М.: Высшая школа, 2000.
22. Вержбицкий В.М. Численные методы. Математический анализ и обыкновенные дифференциальные уравнения. М.: Высшая школа, 2001.

23. **Воеводин В.В.** Вычислительные основы линейной алгебры. М.: Наука, 1977.
24. **Воеводин В.В., Кузнецов Ю.А.** Матрицы и вычисления. М.: Наука, 1984.
25. **Волков Е.А.** Численные методы. М.: Наука, 1987.
26. **Галлягер Р.** Метод конечных элементов. Основы. М.: Мир, 1984.
27. **Гантмахер Ф.Р.** Теория матриц. М.: Наука, 1988.
28. **Гилл Ф., Мюррей У., Райт М.** Практическая оптимизация. М.: Мир, 1985.
29. **Годунов С.К.** Современные аспекты линейной алгебры. Новосибирск: Научная книга, 1997.
30. **Годунов С.К., Рябенький В.С.** Разностные схемы. М.: Наука, 1977.
31. **Голуб Дж., Ван Лоун Ч.** Матричные вычисления. М.: Мир, 1999.
32. **Горинштейн А.М.** Практика решения инженерных задач на ЭВМ. М.: Радио и связь, 1984.
33. **Деккер К., Вервер Я.** Устойчивость методов Рунге—Кутты для жестких нелинейных дифференциальных уравнений. М.: Мир, 1988.
34. **Деклу Ж.** Метод конечных элементов. М.: Мир, 1976.
35. **Демидович Б.П., Марон И.А.** Основы вычислительной математики. СПб.: Лань, 2011.
36. **Демидович Б.П., Марон И.А., Шувалова Э.З.** Численные методы анализа. СПб.: Лань, 2010.
37. **Денисов А.М.** Введение в теорию обратных задач. М.: Изд-во МГУ, 1994.
38. **Джордж А., Лю Дж.** Численное решение больших разреженных систем уравнений. М.: Мир, 1984.
39. **Джоунс У., Трон В.** Непрерывные дроби. М.: Мир, 1985.
40. **Денисис Дж., Шнабель Р.** Численные методы безусловной оптимизации и решения нелинейных уравнений. М.: Мир, 1988.
41. **Жаблон К., Симон Ж.-К.** Применение ЭВМ для численного моделирования в физике. М.: Наука, 1983.
42. **Завьялов Ю.С., Квасов Б.И., Мирошниченко В.Л.** Методы сплайн-функций. М.: Наука, 1980.
43. **Зенкевич О.** Метод конечных элементов в технике. М.: Мир, 1975.
44. **Зенкевич О., Морган К.** Конечные элементы и аппроксимация. М.: Мир, 1986.
45. **Иванов В.В.** Методы вычислений на ЭВМ: Справочное пособие. Киев: Наукова думка, 1986.
46. **Икрамов Х.Д.** Численные методы линейной алгебры. (Решение линейных уравнений). М.: Знание, 1987.
47. **Икрамов Х.Д.** Численные методы для симметричных линейных систем. М.: Наука, 1988.
48. **Икрамов Х.Д.** Вычислительные методы линейной алгебры. (Решение больших разреженных систем уравнений прямыми методами). М.: Знание, 1989.
49. **Икрамов Х.Д.** Несимметричная проблема собственных значений. М.: Наука, 1991.

50. Ильин В.П., Кузнецов Ю.И. Трехдиагональные матрицы и их приложения. М.: Наука, 1985.
51. Калиткин Н.Н. Численные методы. М.: Наука, 1978.
52. Канторович А.В., Акилов Г.Л. Функциональный анализ. М.: Наука, 1977.
53. Карпов В.Я., Корягин Д.А. Пакеты прикладных программ. М.: Знание, 1983.
54. Каханер Д., Моулер К., Нэш С. Численные методы и программное обеспечение. М.: Мир, 1999.
55. Киреев В.И., Пантелеев А.В. Численные методы в примерах и задачах. М.: Изд-во МАИ, 2000.
56. Копченова Н.В., Марон И.А. Вычислительная математика в примерах и задачах. СПб.: Лань, 2009.
57. Косарев В.И. 12 лекций по вычислительной математике. М.: Изд-во МФТИ. Физматкнига, 2000.
58. Краснощеков П.С., Петров А.А. Принципы построения моделей. М.: Изд-во МГУ, 1983.
59. Кронрод А.С. Узлы и веса квадратурных формул. М.: Наука, 1964.
60. Крылов В.И., Бобков В.В., Монастырный П.Н. Вычислительные методы. Т. I. М.: Наука, 1976.
61. Крылов В.И., Бобков В.В., Монастырный П.Н. Вычислительные методы. Т. II. М.: Наука, 1977.
62. Лебедев В.И. Функциональный анализ и вычислительная математика. М.: ВИНИТИ, 1994.
63. Локуциевский О.В., Гавриков М.Б. Начала численного анализа. М.: ТОО «Янус», 1995.
64. Лоусон У., Хенсон Р. Численное решение задач метода наименьших квадратов. М.: Наука, 1986.
65. Люк Ю. Специальные математические функции и их аппроксимации. М.: Мир, 1980.
66. Люстерник Л.А., Червоненкис О.А., Янпольский А.Р. Математический анализ. Вычисление элементарных функций. М.: Физматгиз, 1963.
67. Марчук Г.И. Методы вычислительной математики. СПб.: Лань, 2009.
68. Марчук Г.И., Агошков В.И. Введение в проекционно-сеточные методы. М.: Наука, 1981.
69. Марчук Г.И. Шайдуров В.В. Повышение точности решений разностных схем. М.: Наука, 1979.
70. Митчелл Э., Уэйт Р. Метод конечных элементов для уравнений с частными производными. М.: Мир, 1981.
71. Моисеев Н.Н. Математика ставит эксперимент. М.: Наука, 1979.
72. Морозов В.А. Регулярные методы решения некорректно поставленных задач. М.: Наука, 1987.
73. Ортега Дж., Пул У. Введение в численные методы решения дифференциальных уравнений. М.: Наука, 1986.

74. Ортега Д., Рейнболдт В. Итерационные методы решения нелинейных систем уравнений со многими неизвестными. М.: Мир, 1975.
75. Парлетт Б. Симметричная проблема собственных значений. М.: Мир, 1983.
76. Пирумов У.Г. Численные методы. М. Изд-во МАИ. 1998.
77. Писсанецки С. Технология разреженных матриц. М.: Мир, 1988.
78. Поляк Б.Т. Введение в оптимизацию. М.: Наука, 1983.
79. Попов Ю.П., Самарский А.А. Вычислительный эксперимент. М.: Знание, 1983.
80. Пшеничный Б.П., Данилин Ю.М. Численные методы в экстремальных задачах. М.: Наука, 1975.
81. Райс Дж. Матричные вычисления и математическое обеспечение. М.: Мир, 1984.
82. Ракитский Ю.В., Устинов С.М., Черноруцкий М.Г. Численные методы решения жестких систем. М.: Наука, 1979.
83. Рябенький В.С. Введение в вычислительную математику. М.: Наука, 1994.
84. Самарский А.А. Введение в численные методы. СПб.: Лань, 2009.
85. Самарский А.А. Введение в теорию разностных схем. М.: Наука, 1974.
86. Самарский А.А. Теория разностных схем. М.: Наука, 1977.
87. Самарский А.А., Вабищевич П.Н., Самарская Е.А. Задачи и упражнения по численным методам. М.: Эдиториал УРСС, 2000.
88. Самарский А.А., Гулин А.В. Численные методы. М.: Наука, 1989.
89. Самарский А.А., Николаев Е.С. Методы решения сеточных уравнений. М.: Наука, 1978.
90. Сегерлинг Л. Применение метода конечных элементов. М.: Мир, 1979.
91. Современные численные методы решения обыкновенных дифференциальных уравнений / Под. ред. Дж. Холла и Дж. Уатта. М.: Мир, 1979.
92. Справочник по специальным функциям с формулами, графиками и математическими таблицами / Под ред. М. Абрамовича и И. Стиган. М.: Наука, 1979.
93. Стрэнг Г., Фикс Дж. Теория метода конечных элементов. М.: Мир, 1977.
94. Сухарев А.Г., Тимохов А.В., Федоров В.В. Курс методов оптимизации. М.: Наука, 1986
95. Съярле Ф. Метод конечных элементов для эллиптических задач. М.: Мир, 1980.
96. Тихонов А.Н. Математическая модель. Математическая энциклопедия. Т. 3. М.: Советская энциклопедия, 1982.
97. Тихонов А.Н., Арсенин В.Я. Методы решения некорректных задач. М.: Наука, 1986
98. Тихонов А.Н., Гончарский А.В., Степанов В.В., Ягода Л.Г. Численные методы решения некорректных задач. М.: Наука, 1990.
99. Тихонов А.Н., Костомаров Д.П. Вводные лекции по прикладной математике. М.: Наука, 1984.

100. Тыртышников Е.Е. Краткий курс численного анализа. М.: ВИНИТИ, 1994.
101. Тьюарсон Р. Разреженные матрицы. М.: Мир, 1977.
102. Уилкинсон Дж.Х. Алгебраическая проблема собственных значений. М.: Наука, 1970.
103. Уилкинсон Дж. Х, Райнш К. Справочник алгоритмов на языке Алгол. Линейная алгебра. М.: Машиностроение, 1976.
104. Федоренко Р.П. Введение в вычислительную физику. М.: Изд-во МФТИ, 1994.
105. Флетчер К. Численные методы на основе метода Галеркина. М.: Мир, 1988.
106. Форсайт Дж., Мальcolm M., Моулер К. Машины методы математических вычислений. М.: Мир, 1980.
107. Форсайт Дж., Молер К. Численное решение систем линейных алгебраических уравнений. М.: Мир, 1969.
108. Хайрер Э., Ваннер Г. Решение обыкновенных дифференциальных уравнений. Жесткие и дифференциально-алгебраические задачи. М.: Мир, 1999.
109. Хайрер Э., Нерсетт С., Ваннер Г. Решение обыкновенных дифференциальных уравнений. Нежесткие задачи. М.: Мир, 1990.
110. Хейгман Л., Янг Д. Прикладные итерационные методы. М.: Мир, 1986.
111. Хемминг Р.В. Численные методы для научных работников и инженеров. М.: Наука, 1972.
112. Химмельблау Д. Прикладное нелинейное программирование. М.: Мир, 1975.
113. Чуи К. Введение в вэйвлеты. М.: Мир, 2001.
114. Шул Т. Решение инженерных задач на ЭВМ. М.: Мир, 1982.
115. Эльсгольц Л.Э. Дифференциальные уравнения и вариационное исчисление. М.: Наука, 1969.
116. Эстербю О., Златев З. Прямые методы для разреженных матриц. М.: Мир, 1987.
117. Allen III M.B., Isaacson E.I. Numerical analysis for applied science. JOHN WILEY & SONS INC, 1997.
118. Cheney W., Kincard D. Numerical mathematics and computing. Brooks, Cole Publishing Company, 1999.
119. Gautschi W. Numerical Analysis. An Introduction. Birkhauser, 1997.
120. Kress R. Numerical Analysis. Springer, 1998.
121. Saad Y. Iterative methodes for sparse linear systems. Second Edition. SIAM. Philadelphia, 2000.
122. Stoer J., Bulirsch R. Introduction to Numerical Analysis. Springer-Verlag, New York, 2002.
123. BarretR., Berry M., Chan T.F., Demmel J., Donato J., Dongarra J., Eijkhout V., Pozo R., Romine C., Van der Vorst H. Templates for the Solution of Linear Systems: Building Blocks for Iterative Mettods. SIAM, 1994.

Список литературы, добавленной к 3-му изданию

- 1*. **Арушянин И.О.** Численное решение интегральных уравнений методом квадратур. М.: Изд-во МГУ, 2002.
- 2*. **Баландин М.Ю., Шурина Э.П.** Методы решения СЛАУ большой размерности. Новосибирск: Изд-во НГТУ, 2000.
- 3*. **Васильева А.Б., Тихонов Н.А.** Интегральные уравнения. СПб.: Лань, 2009.
- 4*. **Гавурин М.К.** Лекции по методам вычислений. М.: Наука, 1971.
- 5*. **Деммель Дж.** Вычислительная линейная алгебра. Теория и приложения. М.: Мир, 2001.
- 6*. **Костомаров Д.П., Фаворский А.П.** Вводные лекции по численным методам. М.: Логос, 2004.
- 7*. **Краснов М.Л.** Интегральные уравнения. Введение в теорию. М.: Наука, 1975.
- 8*. **Малла С.** Вэйвлеты в обработке сигналов. М.: Мир, 2005.
- 9*. **Манжиров А.В., Полянин А.Д.** Методы решения интегральных уравнений. Справочник. М.: Факториал, 1999.
- 10*. **Самарский А.А., Гулин А.В.** Численные методы математической физики. М.: Научный мир, 2003.
- 11*. **Самарский А.А., Михайлов В.П.** Математическое моделирование. М.: Физматлит, 2002.
- 12*. **Тыртышников Е.Е.** Матричный анализ и линейная алгебра. М.: Физматлит, 2007.
- 13*. **Фрейзер М.** Введение в вэйвлеты в свете линейной алгебры. М.: БИНОМ. Лаборатория знаний, 2007.
- 14*. **Atkinson K.E.** An Introduction to Numerical Analysis. John Wiley, New York, 1989.
- 15*. **Butcher J.C.** Numerical Methods for Ordinary Differential Equations. Wiley, 2003.
- 16*. **Dahlquist G., Bjørck A.** Numerical Methods in Scientific Computing. Volume 1. SIAM, Philadelphia, 2008.
- 17*. **Dahlquist G., Bjørck A.** Numerical Methods in Scientific Computing. Volume 2. SIAM, Philadelphia, 2008.
- 18*. **Conte S.D., de Boor C.** Elementary numerical analysis. An algorithmic approach. McGraw-Hill Book Company, 1980.
- 19*. **Hackbusch W.** Integral Equations. Theory and Numerical Treatment. Birkhauser Verlag. Basel - Boston - Berlin, 1995.
- 20*. **Heath M.T.** Scientific Computing. An Introductory Survey. McGraw-Hill, Boston, MA, 2002.
- 21*. **Higham N.J.** Accuracy and Stability of Numerical Algorithms. SIAM. Philadelphia. PA, 1996.

- 22* **Hoffman J.D.** Numerical methods for Engineers and Scientists. Marcel Dekker, Inc. New York — Bazel, 2001.
- 23* **Kincaid D., Cheney W.** Numerical Analysis. Brooks/Cole, Pacific Grove, CA, 2002.
- 24* **Quarteroni A., Sacco R., Salery F.** Numerical Mathematics. Springer-Verlag. New York, 2000.
- 25* **Saad Y.** Numerical Methods for Large Eigenvalue Problems. Second Edition, SIAM, Philadelphia, 2003.
- 26* **Van der Vorst H.A.** Iterative Krylov Methods for Large Linear Systems. Cambridge University Press. Cambridge, UK, 2003.
- 27* **Tytychnikov E.** A Brief Introduction to Numerical Analysis. Birkhauser, Boston, 1997.

ПРЕДМЕТНЫЙ УКАЗАТЕЛЬ

- Абсолютная погрешность 27
— вектора 133
— устойчивость 527
- Абстрактный вычислительный алгоритм 69
- Автоматический выбор шага 514
- Адаптивные процедуры численного интегрирования 462
- Алгоритм бинарный 84
— Грама—Шмидта 192
— модифицированный 194
— вычислительно устойчивый 71
— вычислительный 69
— абстрактный 69
— Гира 551
— Краута 200
— неустойчивый 71
— плохо обусловленный 77
— устойчивый 71
— хорошо обусловленный 77
— *LR* 284
— *QR* 279
— основной 279
— со сдвигами 282
- Алгоритмы гибридные 128, 310
— регуляризованные 128, 310
- Анализ ошибок обратный 79
— прямой 78
— статистический 81
- Антиградиент 314
- Антiperеполнение 42
- Апостериорные оценки погрешности 68
- Аппроксимационная теорема
Вейерштрасса 380
- Аппроксимация дифференциального уравнения 490, 567
— — — с p -м порядком 490, 568
— краевых условий 577
- Аппроксимация Паде 421
- Априорные оценки погрешности 67
— — решения краевой задачи 558
— — — разностной схемы 566
- Базис лагранжев 386
— локальный степенной 386
— нормированный степенной 386
— степенной 386
— чебышевский 386
- Базисные многочлены 386
— функции 580
— — — кусочно линейные 586
— — — специальные 589
- Баллистический метод 591
- Безье кривые 425
- Быстрое дискретное преобразование Фурье 398
- Вариационная постановка краевой задачи 579
- Вариационно разностная схема 588
- Вариационные принципы 579
- Вариационный функционал 579
- Вековое уравнение 260
- Вектор Нордсика 553
- Векторы взаимно сопряженные 335
- Верная цифра 29
- Верхняя граница абсолютной погрешности 28
— — относительной погрешности 28
- Вес 467
- Весовая функция 467
- Веса квадратурной формулы 438
— — — Гаусса 454
- Временная постоянная 482
— — локальная 482
- Входное данное 12

- Вычислительная задача** 14, 48
 — корректная 48
 — некорректная 49, 53
 — плохо обусловленная 54
 — хорошо обусловленная 54
 — погрешность 26
Вычислительный алгоритм 69
 — абстрактный 69
 — корректный 69
 — некорректный 70
 — неустойчивый 71
 — устойчивый 71
 — процесс 21
 — эксперимент 22
Вычислительные методы 15, 61
Вэйвлеты 425
- Гибридные алгоритмы** 128, 310
Главный член погрешности
 квадратурной формулы 456
Глобальная погрешность 491
 — полиномиальная интерполяция 380
Градиент 314
Градиентный метод 322
Граница абсолютной погрешности 28
 — относительной погрешности 28
Граничные условия 393
- Данные входные** 12
 — выходные 12
 — исходные 12
Двоичный порядок 39
Двухточечная краевая задача 555
Денормализованные числа 40
Дефект сплайна 390
Диагональное преобладание
 строчное 156
 — столбцовое 206
Дискретизация 63
Дискретная задача Коши 486
 — краевая задача 555, 560
Дискретное преобразование Фурье 396
 — быстрое 398
 — обратное 397
 — прямое 397
- Диссипативная система**
 дифференциальных уравнений 536
Дробление шага 320
Дробь непрерывная 420
 — цепная 420
- Евклидова норма вектора** 132
 — матрицы 135
Естественный кубический сплайн 393
- Жесткая задача Коши** 543
 — система дифференциальных
 уравнений 547
Жорданова форма матрицы 263
- Задача безусловной минимизации** 312
 — вычислительная 14, 48
 — корректная 48
 — некорректная 49
 — плохо обусловленная 54
 — хорошо обусловленная 54
 — дискретной минимизации 345
 — жесткая 543
 — идентификации 13
 — интерполяции 350
 — — обобщенными многочленами 351
 — конечномерная 63
 — линейного программирования 345
 — Коши 478
 — — дискретная 486
 — — для обыкновенного
 дифференциального уравнения
 первого порядка 477
 — — — — — *m*-го порядка 542
 — — — системы обыкновенных
 дифференциальных уравнений
 первого порядка 534
 — — жесткая 543
 — — — для системы дифференциаль-
 ных уравнений 547
 — — краевая двухточечная 555
 — — минимизации оценки погрешности
 интерполяции 365
 — — начальная 476

- Задача наименьших квадратов**
 — линейная 188, 402
 — — — нелинейная 341
 — **нелинейного программирования 345**
 — о наилучшем равномерном приближении 414
 — о понижении степени многочлена 417
 — обратная 13
 — одномерной минимизации 286
 — прямая 12
 — регуляризованная 636
 — условной минимизации 312
Золотое сечение 303
Значащая цифра 28
Интегральная кривая 477
Интегральное тождество 582
Интегральное уравнение Абеля 600
 — — переноса излучения 600
Интегральные уравнения 598
 — Вольтерра второго рода 599
 — — — первого рода 598
 — линейные 599
 — — — нелинейные 599
 — — — Фредгольма второго рода 598
 — — — первого рода 598
Интегрирование дифференциального уравнения 477
Интегро-интерполяционный метод 576
Интервал неопределенности корня нелинейного уравнения 95
 — — точки локального минимума 292
Интервальный анализ 47
Интерполирование 350
Интерполяционный массив 381
 — многочлен 355
 — — — Бесселя 422
 — — — кубический Эрмита 360
 — — — Лагранжа 356
 — — — Ньютона с конечными разностями для интерполяции вперед 379
 — — — — для интерполяции назад 380
 — — — с разделенными разностями 376
 — — обобщенный 351
Интерполяционный многочлен с кратными узлами 359
 — — Эрмита 360
Интерполяция 350
 — билинейная 424
 — глобальная полиномиальная 380
 — квадратичная 357
 — кубическая 357
 — кусочно-полиномиальная 388
 — линейная 357, 424
 — многомерная 423
 — локальная 387
 — обратная 422
 — рациональная 420
 — тригонометрическая 399
Искомое решение 12
Исчезновение порядка 42
Искрепывание 275
Итерационная последовательность 66
 — функция 100
Итерационное уточнение 185
 — — корней нелинейного уравнения 91
Итерационные методы 66
 — двухслойные 218
Итерационный метод одношаговый 92
 — — *k*-шаговый 92
Итерационный метод с оптимальным выбором параметра явный стационарный 220
 — — — — — неявный 224
 — — с чебышевским набором параметров явный 222
 — — — — — неявный 226
 — — — Эйлера 127
 — — процесс 66
Итерация 66
Катастрофическая потеря точности 33, 75
Квазиньютоновские методы 333, 345
Квазиньютоновское условие 334
Квадратурная сумма 438
 — формула 438
 — — — Гаусса 452
 — — — Гаусса—Кронрода 465

- Квадратурная формула
 — интерполяционного типа 447
 — левых прямоугольников 440
 — Ньютона—Котеса 447
 — правых прямоугольников 440
 — прямоугольников элементарная 440
 — составная 440
 — Симпсона элементарная 442
 — составная 442
 — точная для многочленов
 степени m 438
 — трапеций элементарная 441
 — составная 441
 — центральных прямоугольников 440
- Компьютер 6 разрядный десятичный 45
- Конечномерная задача 63
- Конечно разностная схема 521
- Конечно-разностные методы 521
- Конечные разности 367
 — вперед 367
 — назад 372
 — порядка k 367
- Конечные элементы 585
- Константа Лебега 385
- Корень нелинейного уравнения 87
 — простой 87
 — кратный 87
- Корневое условие 524
- Корректность вычислительного алгоритма 69
 — вычислительной задачи 48
- Коэффициент роста 153
- Краевая задача двухточечная 555
 — дискретная 560
- Краевые условия первого рода 556
 — второго рода 556
- Кратность корня нелинейного уравнения 87
 — узла интерполяции 360
- Критерий окончания итерационного процесса 67
- Круги Гершгорина 264
- Кубатурные формулы 474
- Кусочно-полиномиальная интерполяция 388
- Линейная задача наименьших квадратов 188, 402
- Линия уровня 313
- Локализация корней нелинейного уравнения 90
 — собственных значений 264
- Локальная интерполяция 387
 — погрешность 490
- Локальный сплайн 392
- Ломаная Эйлера 497
- Мантисса 39
- Масштабирование 156
- Математическая модель 10
 — гипотетическая 11
 — динамическая 12
 — статическая 12
- Математическое моделирование 10
- Матрица
 — верхняя треугольная 137
 — Гессе 315
 — Гильберта 143
 — Грама 352
 — диагональная 136
 — единичная 136
 — заполненная 138
 — ленточная 139
 — нижняя треугольная 137
 — ортогональная 176
 — отражения 181
 — плотная 138
 — плохо обусловленная 141
 — подобия 262
 — положительно определенная 137
 — простой структуры 263
 — разреженная 138
 — симметричная 137
 — трехдиагональная 139
 — треугольная 137

- Матрица Хаусхолдера** 181
 — Хессенберга 282
 — Якоби 240
Матрицы подобные 262
Машинная бесконечность 42
 — точность 41
Машинное слово 38
 — эпсилон 41
Метод Адамса 517
 — интерполяционный 518
 — экстраполяционный 518
 — Адамса—Башфорта 518
 — Адамса—Моултона 518
 — баланса 576
 — баллистический 591
 — бисекции 97, 307
 — бисекций 284
 — Брайдена 254
 — вращений 176
 — Якоби 284
 — выбора параметра регуляризации по невязке 638
 — Галеркина 583, 619
 — итерированный 621
 — Галлея 127
 — Гаусса 145
 — с выбором главного элемента по всей матрице 154
 — столбцу 152
 — Гаусса—Зейделя 211
Метод Гаусса—Жордана 200
 — Гаусса—Ньютона 342
 — модифицированный 343
 — градиентный 322
 — Давиденко 258
 — Данилевского 261
 — деления отрезка пополам 298
 — деформируемого многогранника 339
 — дифференцирования по параметру 257
 — Дормана—Принса 512
 — замены ядра на вырожденное 628
 — Зейделя 210
 — нелинейный 248
 — Метод золотого сечения 303
 — интегро-интерполяционный 576
 — Канторовича 621
 — итерированный 621
 — касательных 113
 — квадратных корней 168
 — квадратур 611
 — коллокации 626
 — конечных разностей 559
 — элементов 585
 — Крылова 261
 — Ланцюша 285
 — Левенберга—Маркардта 344
 — Леверье 261
 — ложного положения 120, 253
 — минимальных невязок 227
 — минимальных поправок 228
 — Монте—Карло 474
 — Мюллера 128
 — наименьших квадратов 188, 400
 — наискорейшего спуска 228, 323
 — Нумерова 543
 — Ньютона вычисления \sqrt{a} 66
 — минимизации 308, 330
 — модифицированный 251
 — решения нелинейных уравнений 112
 — систем нелинейных уравнений 248
 — упрощенный 119, 252
 — обратной квадратичной интерполяции 127
 — обратных итераций 276
 — с использованием отношения Рэлея 278
 — отражений 181
 — пассивного поиска 296
 — покоординатного спуска 320
 — последовательного исключения неизвестных 145
 — последовательной верхней релаксации 216
 — нижней релаксации 216
 — последовательных замещений 211

- Метод пристрелки 591
 — прогонки 171
 — продолжения по параметру 255
 — простой итерации решения систем линейных алгебраических уравнений 202
 — — — — систем нелинейных уравнений 243
 — — — — нелинейных уравнений 100
 — регуляризации Тихонова 635
 — релаксации 216
 — Ритца 580
 — Ричардсона 222
 — Ромберга 461
 — Рунге—Кутты 506
 — — — неявный m -этапный 514
 — — — четвертого порядка точности 509
 — — — явный m -этапный 508
 — Рунге—Кутты—Фельберга 511
 — секущих 121, 253
 — сеток 559
 — симметричной последовательной верхней релаксации 217
 — сопряженных градиентов 230, 336
 — — — с предобусловливанием 232
 — — — направлений 335
 — Стеффенсена 122, 254
 — степенной 269
 — стрельбы 591
 — установления 258
 — Фибоначчи 301
 — Холецкого 168
 — Хьюна 503
 — Эйлера 487, 496
 — — неявный 488, 530
 — — усовершенствованный 504
 — Эйлера—Коши 503
 — экстраполяции Ричардсона 461
 — Якоби 204
 — SOR 215
 Метода Гаусса ведущий элемент k -го шага 147
 — — главный элемент k -го шага 147
 — — множители k -го шага 147
 Метода Гаусса обратный ход 147
 — — прямой ход 146
 Методы Адамса 517
 — аппроксимации 62
 — вычислительные 15
 — итерационные 66
 — итерирования подпространства 285
 — квазиньютоновские 333, 345
 — конечно-разностные 521
 — линеаризации 63
 — линейные многошаговые 521
 — многозначные 553
 — Монте-Карло 68, 474
 — Нордика 553
 — «овражные» 327
 — переменной метрики 333
 — последовательного поиска 298
 — прогноза и коррекции 504, 520
 — прямого поиска 296, 338
 — прямые 64
 — регуляризации 54
 — Рунге—Кутты 506
 — спуска 318
 — статистических испытаний 68
 — — предиктор-корректор 504
 — точные 65
 — численные 15
 — чисто неявные 550
 — эквивалентных преобразований 61
 Минимальное значение функции 286, 312
 Многочлен интерполяционный 355
 — — кубический Эрмита 360
 — — Ньютона 376, 379
 — — обобщенный 351
 — — с кратными узлами 360
 — Лагранжа 356
 — наилучшего равномерного приближения 414
 — — среднеквадратичного приближения 402
 Многочлен обобщенный 349
 — характеристический 523
 Многочлены наименее уклоняющиеся от нуля 364
 — Чебышева 362

- Модельное уравнение 482
 Множество граничных узлов сетки 560
 — внутренних узлов сетки 560
 — возможных решений 48
 — допустимых входных данных 48
- Надежность программы 84
 Наименьших квадратов задача
 линейная 188, 402
 — — — нелинейная 341, 413
 Наклон сплайна 391
 Направление ньютоновское 330
 — спуска 318
 Начальная задача 476
 Начальное значение 477
 — условие 477
 Начальные значения 486
 — условия 534
 Невязка 131
 Некорректная задача 49, 53
 Нелинейная задача метода наименьших
 квадратов 341
 Непрерывная дробь 420
 Неравенство треугольника 132
 Неустранимая погрешность 26
 Неявный метод Эйлера 488, 530
 Норма вектора 132
 — — евклидова 132
 — матрицы подчиненная 134
 — — евклидова 135
 — Фробениуса 135
 Нормальная система метода
 наименьших квадратов 189, 403
 Нуль-устойчивость 523
- Область абсолютной устойчивости 528
 — неопределенности 79, 243
 — сходимости метода 66
 Обобщенный многочлен 349
 Обратная задача 13
 — интерполяция 422
 — прогонка 173
 Обратный анализ ошибок 79
 — ход метода Гаусса 147
- Обратный ход метода прогонки 173
 Обусловленность вычислительного
 алгоритма 77
 — вычислительной задачи 54
 Односторонние формулы численного
 дифференцирования 432
 Округление 31
 — по дополнению 31
 — усечениям 31
 Определитель Вандермонда 355
 — Грама 353
 Оптимальный пассивный поиск 296
 Остаточный член квадратурной
 формулы 438
 Относительная погрешность 27
 — — вектора 133
 — — точность 28
 Отношение Рэлея 266
 Отражение 181
 Отрезок локализации корня
 нелинейного уравнения 90
 — — — точки локального минимума 287
 — — наблюдения 351
 Оценки погрешности априорные 67
 — — апостериорные 68
 Ошибка 27
 — округления 41
 — представления 41
- Пакет прикладных программ
 проблемно-ориентированный 24
 Пара Гаусса—Кронрода 465
 Параметр пристрелочный 591
 — релаксации 216
 Параметры модели 13
 Переносимость 85
 Переобуславливатель 225
 Переполнение 42
 Пivotирование 152
 Плоское вращение 179
 Плохо обусловленная вычислительная
 задача 54
 — — матрица 141

- Плохо обусловленная система линейных алгебраических уравнений 141
- Плохо обусловленный вычислительный алгоритм 77
- Погрешность 27
- абсолютная 27
 - аппроксимации 62
 - — дискретного уравнения 490
 - — разностного уравнения 567
 - — формулы численного дифференцирования 427
 - — вычислительная 26
 - — квадратурной формулы 438
 - локальная 490
 - метода 26
 - на шаге 490
 - неустранимая 25
 - округления 31
 - относительная 27
 - разностной схемы 568, 569
 - численного метода решения задачи Коши 491
 - — — — глобальная 491
 - — — — локальная 490
- Подобные матрицы 262
- Поверхность уровня 313
- Поддерживаемость 85
- Поле направлений 477
- Полином «движущийся» 387
- устойчивости 528
- Порог машинного нуля 40
- переполнения 40
- Портабельность 85
- Порядок двоичный 39
- сходимости итерационного метода 92
 - точности численного метода решения задачи Коши 491
- Постановка краевой задачи вариационная 579
- — — проекционная 582
- Постоянная временная 482
- Липшица 478
- Правило Гарвика 97
- двойного пересчета 458
- Правило Крамера 83
- Рунге практической оценки погрешности 458, 570
 - трапеций 488
- Предобусловливатель 225, 235
- неявный 236
 - явный 236
- Предобусловливание 225, 236
- Представимое множество компьютера 41
- Преобразование Гивенса 179
- подобия 262
 - Фурье дискретное 397
 - — быстрое 398
 - — обратное 397
 - — — прямое 397
 - Хаусхолдера 181
- Приближенное число 27
- Приведение к виду, удобному для итераций, системы линейных алгебраических уравнений 203
- — — — — нелинейного уравнения 107
- Пример Рунге 382
- Уилкинсона 267
- Принцип максимума 557
- — для системы сеточных уравнений 564
 - — для разностной схемы 564
- Пристрелочные параметры 591
- Проблема «оврагов» 326
- собственных значений полная 260
 - — частичная 260
- Пробная функция 582
- Пробные точки 296
- Прогонка обратная 173
- прямая 172
- Прогоночные коэффициенты 172
- Проекционная постановка краевой задачи 582
- Проекционно-разностная схема 588
- — — специальная 589
- Проекционно-сеточный метод 588
- Проекционные методы 619
- Процесс вычислительный 21
- Либмана 211

- δ^2 — процесс ускорения сходимости 237
 Прямая задача 12
 — прогонка 171
 Прямой анализ ошибок 78
 Прямой ход метода Гаусса 146
 — — — прогонки 172
 Прямые методы 64
 — — — решения проблемы собственных значений 261
 Работоспособность 84
 Разделенные разности 373
 Разложение матрицы на множители — LU 163
 — QR 176, 192
 — SVD 196
 Разности конечные 367
 — — вперед 367
 — — назад 372
 — «против ветра» 575
 — «против потока» 575
 — разделенные 374
 Разностная производная вторая 430
 — — левая 426
 — — правая 426
 — — центральная 427
 — схема 521, 562
 — — однородная 577
 Разностное уравнение 561
 — — k -го порядка линейное однородное с постоянными коэффициентами 523
 Разрядность мантиссы 39
 Рациональная арифметика 47
 Регуляризованные алгоритмы 128, 310
 Рекуррентная формула 66
 Решение краевой задачи для одномерного стационарного уравнения теплопроводности 557
 — — — — — с разрывными коэффициентами 576
 — — — — — — — нелинейного уравнения 87
 Решение обыкновенного дифференциального уравнения первого порядка 477
 Робастность 84
 Сетка 485, 560
 — равномерная 485
 — неравномерная 574
 Сеточные функции 486, 561
 Сдвиги по Рэлею 283
 — — Уилкинсону 283
 Сингулярное разложение матрицы 196
 Сингулярные векторы 198
 — числа 196
 Система линейных алгебраических уравнений плохо обусловленная 141
 — сеточных уравнений 561
 — функций линейно зависимая в точках 351
 — — — независимая в точках 352
 — — — ортогональная на множестве точек 354
 Скалярное произведение векторов 133
 Скорость сходимости итерационного метода квадратичная 92
 — — — — кубическая 92
 — — — — линейная 92
 — — — — сверхлинейная 92
 Собственное значение матрицы 259
 — число матрицы 135, 259
 Собственный вектор матрицы 259
 Спектральный радиус матрицы 209
 Сплайн 390
 — интерполяционный 391
 — кубический 390
 — — естественный 393
 — — фундаментальный 393
 — линейный 390
 — локальный 392
 — степени m 390
 Сплайна дефект 390
 — наклон 391
 Среднеквадратичное уклонение 402

- Стандарт IEEE 40
 Статистический анализ ошибок 81
 Стационарная точка функции 287, 314
 Стационарный итерационный метод 218
 — явный с оптимальным параметром 220
 Степенной метод 269
 — — без сдвигов 269
 — — со сдвигами 275
 Субнормальные числа 40
 Схема Горнера 65
 — единственного деления 146
 — конечно-разностная 521
 — разностная 521, 562
 — полного выбора 154
 Схема частичного выбора 152
 — Эйткена 378
 Сходимость итерационного метода 66, 92
 — — — со скоростью геометрической прогрессии 92
 — метода аппроксимации 62
 — — интерполяции 381
 — к треугольной матрице по форме 280
 — последовательности векторов по направлению 271
 — — — норме 133
 — — — покоординатная 134
 — разностной схемы 569
 — — — с m -м порядком точности 569
 — численного метода решения задачи Коши 491
- Таблица конечных разностей 366
 — разделенных разностей 374
 Теорема Вейерштрасса
 — аппроксимационная 380
 — Гершгорина 265
 — локализации 264
 — сравнения 558, 565
 — Фабера 382
 — Чебышева 415
 Точка минимума глобального 286, 312
 — локального 286, 312
 — строгого локального 286, 312
 — стационарная 287, 314
- Точки пробные 296
 — чебышевского альтернанса 416
 Точность 26, 28
 — абсолютная 28
 — двойная 44
 — — расширенная 44
 — машинная 41
 — относительная 28
 — компьютера относительная 41
 Точные методы 65
 Траектория спуска 319
- Угловой коэффициент 477
 Узлы интерполяции 350
 — — кратные 360
 — квадратурной формулы 438
 — — — Гаусса 454
 — равноотстоящие 366
 — сетки 560
 — — внутренние 560
 — — граничные 560
 Уклонение многочлена от нуля 363
 — среднеквадратичное 402
 Унимодальная функция 289
 Упрощенный метод Ньютона 119, 252
 Уравнение диффузии одномерное 556
 — линейное однородное разностное с постоянными коэффициентами 523
 — модельное 482
 — теплопроводности одномерное стационарное 555
 — характеристическое 523
 — Эйлера 579
 Уравновешивание 156
 Усечение 31
 Условие диагонального преобладания 156, 206
 — корневое 524
 — Липшица 478, 535
 — — одностороннее 479, 536
 — «отсутствия узла» 393
 Условия граничные 393
 — краевые второго рода 556
 — — первого рода 556

- Усовершенствованный метод Эйлера 504
 Устойчивость алгоритма
 — вычислительная 71
 — по входным данным 70
 — разностной схемы 567
 — решения вычислительной задачи 50
 — — — — абсолютная 53
 — — — — относительная 53
 — — — — по входным данным 50
 — решения задачи Коши
 асимптотическая 484
 — — — на конечном отрезке по
 начальным значениям 481
 — — — — — и правой
 части 483
 — — — — по Ляпунову 484
 — численного метода решения задачи
 Коши 488
 — — — — абсолютная 528
- Форма матрицы жорданова 263
 — Хессенберга 281
- Формула Ньютона—Лейбница 437
 — парабол 442
 — рекуррентная 66
 — Шермана—Моррисона 201
 — Эрмита 470
- Формулы дифференцирования назад 551
 — квадратурные 438
 — кубатурные 474
 — численного дифференцирования 426
 — — — односторонние 432
 — Филона 471
- Фредгольма теория 604
- Фундаментальный кубический
 сплайн 393
- Функция весовая 467
 — — Лагерра 467
 — — Эрмита 467
 — — Якоби 467
 — выпуклая 315
 — овражная 327
 — пробная 582
 — сеточная 485, 561
- Функция сильно выпуклая 315
 — — строго выпуклая 315
 — — унимодальная 289
 — — целевая 286, 312
- Характеристический многочлен 523
- Характеристическое уравнение 260, 523
- Хорошо обусловленная вычислительная
 задача 54
- Хорошо обусловленный
 вычислительный алгоритм 77
- Целевая функция 286, 312
- Цепная дробь 420
- Цифра значащая 28
 — — верная 29
- Числа Фибоначчи 301
- Численные методы 15
- Численный метод решения
 задачи Коши 485
 — — — — абсолютнно
 устойчивый 528
 — — — — A -устойчивый 530
 — — — — $A(\alpha)$ -устойчивый 550
 — — — — для систем
 дифференциальных уравнений
 первого порядка 540
 — — — — k -шаговый 486
 — — — — линейный
 многошаговый 521
 — — — — многошаговый 487
 — — — — неявный 488
 — — — — нуль-устойчивый 523
 — — — — одношаговый 487
 — — — — самостартующий 487
 — — — — сходящийся 491
 — — — — устойчивый 489
 — — — — на конечном
 отрезке 489
 — — — — явный 488
- Число денормализованное 40
 — жесткости 547
 — нормализованное 39

- Число обусловленности**
— вычислительной задачи 54
— — — абсолютное 55
— — — относительное 55
— — вычислительного алгоритма 77
— — задачи вычисления многочлена с приближенно заданными коэффициентами 387
— — матрицы 141
— — естественное 140
— — стандартное 141
— приближенное 27
— субнормальное 40
— NaN 45
- Шаг сетки** 485
— спуска 319
— таблицы 366
- Шаг формулы численного дифференцирования** 426
- Ширина ленты** 139
- Экономичность алгоритма** 82
- Экономизация степенных рядов** 418
- Экстраполяция** 351
— Ричардсона 461
- Явный итерационный метод**
с чебышевским набором параметров 222
- Явный стационарный итерационный метод** с оптимальным параметром 220
- Ядро интегрального уравнения** 599
— — — вырожденное 627

ОГЛАВЛЕНИЕ

Предисловие к первому изданию	3
Предисловие ко второму изданию	7
Предисловие к третьему изданию	8
Глава 1. Математическое моделирование и решение прикладных задач с применением компьютера	9
1.1. Математическое моделирование и процесс создания математической модели	10
1.2. Основные этапы решения прикладной задачи с применением компьютера	17
1.3. Вычислительный эксперимент	22
1.4. Дополнительные замечания	24
Глава 2. Введение в элементарную теорию погрешностей	25
2.1. Источники и классификация погрешностей результата численного решения задачи	25
2.2. Приближенные числа. Абсолютная и относительная погрешности	26
2.3. Погрешность арифметических операций над приближенными числами	32
2.4. Погрешность функции	34
2.5. Особенности машинной арифметики	37
2.6. Дополнительные замечания	46
Глава 3. Вычислительные задачи, методы и алгоритмы.	
Основные понятия	48
3.1. Корректность вычислительной задачи	48
3.2. Обусловленность вычислительной задачи	54
3.3. Вычислительные методы	61
3.4. Корректность вычислительных алгоритмов	69
3.5. Чувствительность вычислительных алгоритмов к ошибкам округления	73
3.6. Различные подходы к анализу ошибок	78
3.7. Требования, предъявляемые к вычислительным алгоритмам	82
3.8. Дополнительные замечания	86
Глава 4. Методы отыскания решений нелинейных уравнений	87
4.1. Постановка задачи. Основные этапы решения	87
4.2. Обусловленность задачи вычисления корня	93
4.3. Метод бисекции	97
4.4. Метод простой итерации	100
4.5. Обусловленность метода простой итерации	108
4.6. Метод Ньютона	112
4.7. Модификации метода Ньютона	118
4.8. Дополнительные замечания	127

Гла́ва 5. Пря́мые ме́тоды реше́ния систе́м лине́йных алгебра́ических уравнений	130
5.1. Постановка задачи	130
5.2. Нормы вектора и матрицы	131
5.3. Типы используемых матриц	136
5.4. Обусловленность задачи решения системы линейных алгебраических уравнений	139
5.5. Метод Гаусса	145
5.6. Метод Гаусса и решение систем уравнений с несколькими правыми частями, обращение матриц, вычисление определителей	157
5.7. Метод Гаусса и разложение матрицы на множители. LU-разложение .	160
5.8. Метод Холецкого (метод квадратных корней)	168
5.9. Метод прогонки	171
5.10. QR-разложение матрицы. Методы вращений и отражений	176
5.11. Итерационное уточнение	185
5.12. Линейная задача метода наименьших квадратов	188
5.13. Дополнительные замечания	199
Гла́ва 6. Итера́ционные ме́тоды реше́ния систе́м лине́йных алгебра́ических уравнений	202
6.1. Метод простой итерации	202
6.2. Метод Зейделя	210
6.3. Метод последовательной верхней релаксации	215
6.4. Другие двухслойные итерационные методы	218
6.5. Метод сопряженных градиентов	230
6.6. Дополнительные замечания	234
Гла́ва 7. Методы оты́скания реше́ний систе́м нелине́йных уравнений	238
7.1. Постановка задачи. Основные этапы решения	238
7.2. Метод простой итерации	243
7.3. Метод Ньютона для решения систем нелинейных уравнений	248
7.4. Модификации метода Ньютона	251
7.5. О некоторых подходах к решению задач локализации и отыскания решений систем нелинейных уравнений	254
7.6. Дополнительные замечания	258
Гла́ва 8. Методы реше́ния пробле́мы собственных значе́ний	259
8.1. Постановка задачи. Некоторые вспомогательные сведения	259
8.2. Степенной метод	269
8.3. Метод обратных итераций	275
8.4. QR-алгоритм	279
8.5. Дополнительные замечания	284
Гла́ва 9. Методы одномерной минимиза́ции	286
9.1. Задача одномерной минимизации	286
9.2. Обусловленность задачи минимизации	292

9.3. Методы прямого поиска. Оптимальный пассивный поиск. Метод деления отрезка пополам. Методы Фибоначчи и золотого сечения	296
9.4. Метод Ньютона и другие методы минимизации гладких функций	307
9.5. Дополнительные замечания	311
Глава 10. Методы многомерной минимизации	312
10.1. Задача безусловной минимизации функции многих переменных	312
10.2. Понятие о методах спуска. Покоординатный спуск	318
10.3. Градиентный метод	322
10.4. Метод Ньютона	330
10.5. Метод сопряженных градиентов	335
10.6. Методы минимизации без вычисления производных	338
10.7. Нелинейная задача метода наименьших квадратов	341
10.8. Дополнительные замечания	345
Глава 11. Приближение функций и смежные вопросы	347
11.1. Постановка задачи приближения функций	347
11.2. Интерполяция обобщенными многочленами	350
11.3. Полиномиальная интерполяция. Многочлен Лагранжа	355
11.4. Погрешность интерполяции	357
11.5. Интерполяция с кратными узлами	359
11.6. Минимизация оценки погрешности интерполяции. Многочлены Чебышева	361
11.7. Конечные разности	366
11.8. Разделенные разности	373
11.9. Интерполяционный многочлен Ньютона. Схема Эйткена	376
11.10. Обсуждение глобальной полиномиальной интерполяции. Понятие о кусочно-полиномиальной интерполяции	380
11.11. Интерполяция сплайнами	389
11.12. Понятие о дискретном преобразовании Фурье и тригонометрической интерполяции	396
11.13. Метод наименьших квадратов	400
11.14. Равномерное приближение функций	414
11.15. Дробно-рациональные аппроксимации и вычисление элементарных функций	419
11.16. Дополнительные сведения об интерполяции	422
11.17. Дополнительные замечания	425
Глава 12. Численное дифференцирование	426
12.1. Простейшие формулы численного дифференцирования	426
12.2. О выводе формул численного дифференцирования	431
12.3. Обусловленность формул численного дифференцирования	433
12.4. Дополнительные замечания	436
Глава 13. Численное интегрирование	437
13.1. Простейшие квадратурные формулы	437
13.2. Квадратурные формулы интерполяционного типа	446

13.3. Квадратурные формулы Гаусса	452
13.4. Апостериорные оценки погрешности. Понятие об адаптивных процедурах численного интегрирования	455
13.5. Вычисление интегралов в нерегулярных случаях	465
13.6. Дополнительные замечания	473
Глава 14. Численные методы решения задачи Коши для обыкновенных дифференциальных уравнений	476
14.1. Задача Коши для дифференциального уравнения первого порядка	477
14.2. Численные методы решения задачи Коши. Основные понятия и определения	485
14.3. Использование формулы Тейлора	494
14.4. Метод Эйлера	496
14.5. Модификации метода Эйлера второго порядка точности	501
14.6. Методы Рунге—Кутты	506
14.7. Линейные многошаговые методы. Методы Адамса	516
14.8. Устойчивость численных методов решения задачи Коши	522
14.9. Неявный метод Эйлера	530
14.10. Решение задачи Коши для систем обыкновенных дифференциальных уравнений и дифференциальных уравнений m -го порядка	534
14.11. Жесткие задачи	543
14.12. Дополнительные замечания	552
Глава 15. Решение двухточечных краевых задач	555
15.1. Краевые задачи для одномерного стационарного уравнения теплопроводности	555
15.2. Метод конечных разностей: основные понятия	559
15.3. Метод конечных разностей: аппроксимации специального вида	573
15.4. Понятие о проекционных и проекционно-разностных методах. Методы Ритца и Галеркина. Метод конечных элементов	579
15.5. Метод пристрелки	591
15.6. Дополнительные замечания	596
Глава 16. Численное решение интегральных уравнений	598
16.1. Понятие об интегральных уравнениях	598
16.2. Примеры задач, приводящих к интегральным уравнениям	600
16.3. Интегральные уравнения Фредгольма второго рода. Введение в теорию	603
16.4. Интегральные уравнения Фредгольма второго рода. Метод квадратур	611
16.5. Интегральные уравнения Фредгольма второго рода. Проекционные методы	619
16.6. Интегральные уравнения Фредгольма второго рода. Методы наименьших квадратов и коллокации	625
16.7. Интегральные уравнения Фредгольма второго рода. Метод замены ядра на вырожденное	627
16.8. Интегральные уравнения Вольтерра второго рода	629
16.9. Интегральные уравнения Фредгольма первого рода	635

16.10. Интегральные уравнения Вольтерра первого рода.	639
16.11. Нелинейные интегральные уравнения	642
16.12. Дополнительные замечания.	644
 Список литературы.	648
Предметный указатель	655

*Андрей Авенирович АМОСОВ,
Юлий Андреевич ДУБИНСКИЙ,
Наталья Васильевна КОПЧЕНОВА*

ВЫЧИСЛИТЕЛЬНЫЕ МЕТОДЫ

*Учебное пособие
Издание четвертое, стереотипное*

Выпускающие Т. С. Симонова, Н. В. Черезова

ЛР № 065466 от 21.10.97
Гигиенический сертификат 78.01.07.953.П.007216.04.10
от 21.04.2010 г., выдан ЦГСЭН в СПб

Издательство «ЛАНЬ»
lan@lanbook.ru; www.lanbook.com
192029, Санкт-Петербург, Общественный пер., 5.
Тел./факс: (812) 412-29-35, 412-05-97, 412-92-72.
Бесплатный звонок по России: 8-800-700-40-71

ГДЕ КУПИТЬ

ДЛЯ ОРГАНИЗАЦИЙ:

Для того, чтобы заказать необходимые Вам книги, достаточно обратиться в любую из торговых компаний Издательского Дома «ЛАНЬ»:

по России и зарубежью
•ЛАНЬ-ТРЕЙД•, 192029, Санкт-Петербург, ул. Крупской, 13
тел.: (812) 412-85-78, 412-14-45, 412-85-82; тел./факс: (812) 412-54-93
e-mail: trade@lanbook.ru; ICQ: 446-869-967
www.lanpb1.spb.ru/price.htm

в Москве и в Московской области
•ЛАНЬ-ПРЕСС•, 109263, Москва, 7-я ул. Текстильщиков, д. 6/19
тел.: (499) 178-65-85; e-mail: lanpress@lanbook.ru

в Краснодаре и в Краснодарском крае
•ЛАНЬ-ЮГ•, 350072, Краснодар, ул. Жлобы, д. 1/1
тел.: (861) 274-10-35; e-mail: lankrd98@mail.ru

ДЛЯ РОЗНИЧНЫХ ПОКУПАТЕЛЕЙ:

интернет-магазины:

Издательство «Лань»: <http://www.lanbook.com>
•Сова•: <http://www.symplex.ru>; •Озон.ru•: <http://www.ozon.ru>
•Библион•: <http://www.biblion.ru>

Подписано в печать 14.11.13.
Бумага офсетная. Гарнитура Школьная. Формат 84×108 $\frac{1}{32}$.
Печать офсетная. Усл. п. л. 35,28. Тираж 1000 экз.

Заказ №2404.

**Отпечатано в полном соответствии
с качеством предоставленных диапозитивов**
в ОАО «Издательско-полиграфическое предприятие «Правда Севера».
163002, г. Архангельск, пр. Новгородский, д. 32.
Тел./факс (8182) 64-14-54; www.ippps.ru

А. А. Амосов
Ю. А. Дубинский
Н. В. Колченова

ВЫЧИСЛИТЕЛЬНЫЕ МЕТОДЫ



Издательство «Лань»
победитель конкурса по качеству
«Сделано в Санкт-Петербурге»

ISBN 978-5-8114-1623-3



9 785811 416233

